

## 1 Supervised Learning

### 1.1 Definitions

Given  $\mathcal{D} := \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$ ,  $y^{(i)} = f(\mathbf{x}^{(i)}) + \epsilon$ , find  $\hat{f}$  s.t.  $\hat{f}(\mathbf{x}^{(i)}) \approx f^{(i)}(\mathbf{x}^{(i)})$  by  $\hat{f} = \arg \min_{h \in \mathcal{H}} \mathcal{L}(h)$  where  $\mathcal{H}$  is parameterized by  $\mathbf{w}$ .

Regression losses:  $l(y, y') =$

$$\text{Squared } (l_2): (y - y')^2 \quad \text{Huber: } \begin{cases} \frac{1}{2}(y - y')^2 & |y - y'| \leq \delta \\ \delta(|y - y'| - \frac{1}{2}\delta) & \text{otherwise} \end{cases}$$

$$\text{Absolute } (l_1): |y - y'| \quad \epsilon\text{-insensitive: } \begin{cases} 0 & |y - y'| \leq \epsilon \\ |y - y'| - \epsilon & \text{otherwise} \end{cases}$$

Classification losses:  $l(y, y') =$

$$\text{Zero-one: } \begin{cases} 1 & y \neq y' \\ 0 & \text{otherwise} \end{cases} \quad \text{Hinge: } \max(0, 1 - yy')$$

Expected loss consists of bias (average difference from best predictor), variance (sensitivity to choice of  $\mathcal{D}$ ) and Bayes error (irreducible). **Bayes optimal predictors** achieve Bayes error (impossible in practice).

$$\mathbb{E}_{\mathcal{D}} \mathbb{E}_y [l(h(\mathbf{x}; \mathcal{D}) - y)^2 | \mathbf{x}] = \underbrace{\mathbb{E}_{\mathcal{D}} [h - h_*]^2}_{\text{bias}} + \underbrace{\text{Var}(h)}_{\text{variance}} + \underbrace{\text{Var}(y)}_{\text{Bayes error}}$$

$$\mathcal{L}_{\text{gen}}(h) = \mathbb{E}_{(\mathbf{x}, y) \sim \text{Pr}(\mathbf{x}, y)} [l(h(\mathbf{x}), y)] \quad \mathcal{L}_{\text{emp}}(h) = \frac{1}{N} \sum_{i=1}^N l(h(\mathbf{x}^{(i)}), y^{(i)})$$

We want to minimize  $\mathcal{L}_{\text{gen}}(h)$  but we can't because we don't know  $\text{Pr}(\mathbf{x}, y)$ . **Empirical risk minimization** minimizes  $\mathcal{L}_{\text{emp}}(h)$  instead.  $\mathcal{L}_{\text{emp}}(h) \rightarrow \mathcal{L}_{\text{gen}}(h)$  as  $N \rightarrow \infty$  due to WLLN.

**$\nu$ -fold cross validation:** Split  $\mathcal{D}$  into  $\nu$  equal folds, for each fold train on all others and evaluate on current. Increasing  $\nu$  lowers bias but increases variance of error estimate.  $\nu = N$  is **leave-one-out cross validation**.

### 1.2 $k$ -Nearest Neighbours

$$\hat{f}(\mathbf{x}^*) = \frac{1}{k} \sum_{i \in N_k(\mathbf{x}^*)} y^{(i)} \quad \text{or} \quad \frac{1}{\sum_i w_i} \sum_{i \in N_k(\mathbf{x}^*)} w_i y^{(i)}, \quad w_i = \frac{1}{\text{dist}(\mathbf{x}^*, \mathbf{x}^{(i)})}$$

where  $N_k(\mathbf{x}^*)$  are  $k$  training samples closest to  $\mathbf{x}^*$ . Use most common label for classification, decrease  $k$  if tied or use odd  $k$ . Smaller  $k$  gives more complex model. Rule of thumb:  $k \leq \sqrt{N}$ . **Normalize**  $x_i \leftarrow (x_i - \mu_i) / \sigma_i$  per-dimension if scale doesn't matter. Distance metrics:  $\text{dist}(\mathbf{x}, \mathbf{z}) =$

$$\text{Minkowski: } \left( \sum_{i=1}^D |x_i - z_i|^p \right)^{\frac{1}{p}} \quad \text{Mahalanobis: } \sqrt{(\mathbf{x} - \mathbf{z})^T \Sigma^{-1} (\mathbf{x} - \mathbf{z})}$$

**Complexity:**  $\mathcal{O}(ND + N \log N)$  or  $\mathcal{O}(D \log N)$  ( $k$ -d tree,  $D \ll N$ ).

**Curse of dimensionality:** Number of points needed to keep neighbour distance constant grows exponentially with  $D$ . (Use dimensionality reduction.) Result can be made **probabilistic** by considering distribution of labels. Apply Laplace smoothing to avoid zero probabilities.

### 1.3 Linear Regression

$$\hat{f}(\mathbf{x}, \mathbf{w}) = \mathbf{w}_0 + \sum_{j=1}^D w_j x_d = \mathbf{w}^T \mathbf{x}$$

where  $x_0 = 1$ ,  $\mathbf{x} = \{x_0, \dots, x_D\}^T$ ,  $\mathbf{X}_j = x_j^{(i)}$ ,  $\mathbf{y} = \{y^{(1)}, \dots, y^{(N)}\}^T$  then

$$\hat{\mathbf{y}} = \mathbf{X} \mathbf{w} \quad \hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{D+1}} \|\mathbf{y} - \mathbf{X} \mathbf{w}\|_2^2 \implies \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} = \mathbf{X}^T \mathbf{y}$$

Requires  $\mathcal{O}(ND) + \mathcal{O}(N)$  memory to solve (additional  $\mathcal{O}(ND)$  for some methods). Compute  $\mathbf{X}^T \mathbf{X}$  and  $\mathbf{X}^T \mathbf{y}$  term-by-term for large problems.

**Overdetermined problems:**  $N > D + 1$ ,  $\mathbf{X}$  is tall. Solve for  $\mathbf{w}$  by: **Cholesky:**  $\mathbf{X}^T \mathbf{X} = \mathbf{R}^T \mathbf{R}$  where  $\mathbf{R} \in \mathbb{R}^{(D+1) \times (D+1)}$  upper-triangular. Then  $\mathbf{z} = \mathbf{R}^{-T} \mathbf{X}^T \mathbf{y}$  and  $\mathbf{w} = \mathbf{R}^{-1} \mathbf{z}$  by back substitution. Add  $\lambda \mathbf{1}$  to  $\mathbf{X}^T \mathbf{X}$  to avoid singularity ( $l_2$  regularization). Fastest but squares condition number. Complexity:  $\mathcal{O}(N(D+1)^2 + \frac{1}{3}(D+1)^3)$

**Economic QR:**  $\mathbf{X} = \mathbf{Q} \mathbf{R}$  where  $\mathbf{Q} \in \mathbb{R}^{N \times (D+1)}$  is orthonormal,  $\mathbf{R}$  as above. Then  $\mathbf{w} = \mathbf{R}^{-1} \mathbf{X}^T \mathbf{y}$ . Can fail when  $\mathbf{X}$  is nearly rank-deficient. Complexity:  $\mathcal{O}(2N(D+1)^2 + \frac{2}{3}(D+1)^3)$

**SVD:**  $\mathbf{X} = \mathbf{U} \Sigma \mathbf{V}^T = \mathbf{U}_1 \Sigma_1 \mathbf{V}_1^T$  where  $\mathbf{U} \in \mathbb{R}^{N \times N}$ ,  $\mathbf{V} \in \mathbb{R}^{(D+1) \times (D+1)}$  are orthogonal and  $\Sigma \in \mathbb{R}^{N \times (D+1)}$  is diagonal. Then  $\mathbf{w} = \mathbf{V} \Sigma_1^{-1} \mathbf{U}_1^T \mathbf{y}$

$\Sigma_1^{-1}$  is  $\sum_{i=1}^{D+1} \frac{\mathbf{u}_i \mathbf{u}_i^T}{\sigma_i}$  where  $\mathbf{u}_i, \mathbf{v}_i$  denotes column  $i$ . Truncate for rank-deficient  $\mathbf{X}$ .

Slowest and most stable. Complexity:  $\mathcal{O}(2N(D+1)^2 + 11N(D+1)^3)$

**Pseudoinverse:**  $\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{V} \Sigma^{-1} \mathbf{U}^T$  then  $\mathbf{w} = \mathbf{X}^\dagger \mathbf{y}$ .

**Condition number:**  $\kappa(\mathbf{X}) = \|\mathbf{X}\| \|\mathbf{X}^\dagger\| = \frac{\sigma_{\max}}{\sigma_{\min}}$ . Computing  $\mathbf{X}^T \mathbf{X}$  squares it. Adding  $\lambda \mathbf{1}$  lowers it. 1 digit of precision is lost per power of 10 in  $\kappa$ . Cholesky squares  $\kappa$ .

**Underdetermined problems:**  $N < D + 1$ ,  $\mathbf{X}$  is wide. Instead solve  $\hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{D+1}} \|\mathbf{w}\|_2^2$  subject to  $\mathbf{X} \mathbf{w} = \mathbf{y} \implies \mathbf{w} = \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{y}$ .

For QR,  $\mathbf{X}^T = \mathbf{Q} \mathbf{R}$  and  $\mathbf{w} = \mathbf{Q} \mathbf{R}^{-T} \mathbf{y}$ ; for SVD,  $\mathbf{X}^T = \mathbf{U}_1 \Sigma_1 \mathbf{V}_1^T$  and  $\mathbf{w} = \mathbf{U}_1 \Sigma_1^{-1} \mathbf{V}_1^T \mathbf{y}$ .

For classification, use model  $\hat{y} = \sigma(\mathbf{w}^T \mathbf{x})$  where  $\sigma(z) = \frac{1}{1 + e^{-z}}$ . Use cross-entropy instead of squared loss (logistic regression).

$$\mathcal{L}_{CE}(\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \left[ -y^{(i)} \log \hat{y}^{(i)} - (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}) \right]$$

No closed-form solution exists. Gradient is  $\sum_{i=1}^N (y^{(i)} - \sigma(\mathbf{w}^T \mathbf{x}^{(i)})) \mathbf{x}^{(i)}$ .

### 1.4 Generalized Linear Models (GLMs)

$$\hat{f}(\mathbf{x}, \mathbf{w}) = \mathbf{w}_0 + \sum_{i=1}^{M-1} w_i \phi_i(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

where  $\phi(\mathbf{x}) = 1$ ,  $\phi(\mathbf{x}) = \{\phi_0(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x})\}^T$ ,  $\Phi_j = \phi_j(\mathbf{x}^{(i)})$  then

$$\hat{\mathbf{y}} = \Phi \mathbf{w} \quad \hat{\mathbf{w}} = \arg \min_{\mathbf{w} \in \mathbb{R}^{D+1}} \|\mathbf{y} - \Phi \mathbf{w}\|_2^2 \implies \Phi^T \Phi \hat{\mathbf{w}} = \Phi^T \mathbf{y}$$

Large  $M$  leads to overfitting. Can be counteracted by **regularization:**

$$\mathbf{w} = \arg \min_{\mathbf{w} \in \mathbb{R}^M} \|\mathbf{y} - \Phi \mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \implies (\Phi^T \Phi + \lambda \mathbf{1}) \mathbf{w} = \Phi^T \mathbf{y}$$

SVD:  $\Phi = \mathbf{U} \Sigma \mathbf{V}^T \implies \mathbf{w}(\lambda) = \mathbf{V} (\Sigma^2 + \lambda \mathbf{1})^{-1} \Sigma^T \mathbf{U}^T \mathbf{y} = \sum_{i=1}^M v_i \frac{\sigma_i \mathbf{u}_i^T \mathbf{y}}{\sigma_i^2 + \lambda}$

### 1.5 Dual Representation of GLMs

$$\hat{f}(\mathbf{x}, \boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i k(\mathbf{x}, \mathbf{x}^{(i)}) = \mathbf{k}^T(\mathbf{x}) \boldsymbol{\alpha}, \quad \boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{1})^{-1} \mathbf{y}$$

where  $\mathbf{k}(\mathbf{x}) = \{k(\mathbf{x}^{(1)}, \mathbf{x}), \dots, k(\mathbf{x}^{(N)}, \mathbf{x})\}^T$ ,  $K_{ij} = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ . Relation to primal:  $k(\mathbf{x}, \mathbf{z}) = \phi^T(\mathbf{x}) \phi(\mathbf{z})$ ,  $\mathbf{w} = \Phi^T \boldsymbol{\alpha}$ ,  $\alpha_i = -\frac{1}{\lambda} (\mathbf{w}^T \phi(\mathbf{x}^{(i)}) - y^{(i)})$ . Setting  $\lambda = 0$  leads to an interpolation model, i.e. it reproduces training targets given training inputs.

Kernels define an inner product (measure of distance) over Hilbert space  $\mathcal{F}$ :

$$k(\mathbf{x}, \mathbf{z}) = \langle \mathbf{u}(\mathbf{x}), \mathbf{u}(\mathbf{z}) \rangle = \iint \mathbf{u}(\mathbf{x}) k(\mathbf{x}, \mathbf{z}) \mathbf{u}(\mathbf{z}) d\mathbf{x} d\mathbf{z} \geq 0$$

$$k(\mathbf{x}, \mathbf{x}) \geq 0 \quad k(\mathbf{z}, \mathbf{z}) \leq \sqrt{k(\mathbf{x}, \mathbf{x}) k(\mathbf{z}, \mathbf{z})}$$

Example kernels:  $r = \|\mathbf{x} - \mathbf{z}\|$ ,  $k(\mathbf{x}, \mathbf{z}) =$

$$\text{Linear: } \mathbf{x}^T \mathbf{z} \quad \text{Polynomial: } (1 + \mathbf{x}^T \mathbf{z})^n$$

$$\text{Iso. Gaussian: } e^{-\frac{1}{2} r^2} \quad \text{Aniso. Gaussian: } e^{-(\mathbf{x} - \mathbf{z})^T \boldsymbol{\Theta}^{-1} (\mathbf{x} - \mathbf{z})}$$

$$\text{Multiquadratic: } \sqrt{1 + \frac{r^2}{\theta}} \quad \text{Inv. MQ: } \frac{1}{\sqrt{1 + \frac{r^2}{\theta}}}$$

$$\text{Matern } C^0: \exp\left(-\frac{r}{\theta}\right) \quad \text{Matern } C^2: \frac{1}{1 + \frac{r}{\theta}} \exp\left(-\frac{r}{\theta}\right)$$

$$\text{Matern } C^4: \left(3 + 3\frac{r}{\theta} + \left(\frac{r}{\theta}\right)^2\right) \exp\left(-\frac{r}{\theta}\right)$$

All **radial basis function (RBF)** kernels depend only on  $r$  and produce positive definite  $\mathbf{K}$  (except multiquadratics).  $\mathbf{K}$  is non-singular if data points are distinct.

Kernel selection is informed by prior knowledge of the target, e.g. degree of smoothness, periodicity, other function properties. Use kernel metrics when  $M \gg N$  and explicit basis functions when  $N \gg M$ ; if both are high use sparse models (e.g. orthogonal marching pursuit). Time and memory complexity of kernel GLMs scale as  $\mathcal{O}(N^3)$ .

Any algorithm can be kernelized by mapping from input to feature space via a basis  $\phi: \mathcal{X} \rightarrow \mathcal{F}$ , then replacing any inner product  $\phi^T(\mathbf{x}) \phi(\mathbf{z})$  with the kernel  $k(\mathbf{x}, \mathbf{z})$ .

### 1.6 Sparse GLMs

**Orthogonal marching pursuit:** Greedy sparse regression. Pick basis functions one at a time using metric  $J(\phi_i) = \frac{(\Phi^T \mathbf{r}^{(k)})^2}{\Phi_i^T \Phi_i}$  where  $\Phi_i$  denotes column  $i$ . Solve  $\phi^{(k)}(\mathbf{w}^{(k)}) \approx \mathbf{y}$  for updated weights and compute residual  $\mathbf{r}^{(k)} = \mathbf{y} - \Phi^{(k)}(\mathbf{w}^{(k)})$ .

**Leave-one-out-error:** Let  $\mathbf{A} = \mathbf{K}(\mathbf{K} + \lambda \mathbf{1})^{-1} = \Phi^T(\Phi \Phi^T + \lambda \mathbf{1})^{-1} \Phi$  then

$$y^{(i)} - \hat{f}^{(i)}(\mathbf{x}^{(i)}) = \frac{y^{(i)} - \hat{f}(\mathbf{x}^{(i)})}{1 - A_{ii}} \implies \text{LOO} = \frac{1}{N} \sum_{i=1}^N \left( \frac{y^{(i)} - \hat{f}(\mathbf{x}^{(i)})}{1 - A_{ii}} \right)^2$$

for least-squares error with  $l_2$  regularization.

$l_1$  regularization typically gives sparser models, but no analytic solution.

### 1.7 Parameter/Density Estimation

Given  $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$  find its generating distribution, assuming it is a member of  $\mathcal{P}_{\boldsymbol{\theta}} = \{p(\mathbf{x} | \boldsymbol{\theta}) | \boldsymbol{\theta} \in \Gamma\}$  parameterized by  $\boldsymbol{\theta}$ .

**Maximum likelihood:** Assuming no prior, i.i.d. data:

$$\hat{\boldsymbol{\theta}}_{\text{ML}} = \arg \max_{\boldsymbol{\theta} \in \Gamma} p(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)} | \boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta} \in \Gamma} \sum_{i=1}^N \log p(\mathbf{x}^{(i)} | \boldsymbol{\theta})$$

Assuming additive Gaussian noise on  $f(\mathbf{x}^{(i)}, \mathbf{w})$  is equivalent to regression with  $l_2$  loss; Laplacian noise is  $l_1$  loss. However resulting variance is constant, which is not reasonable (should be smaller where there are more data points). **Maximum a posteriori:** Assuming prior  $p(\boldsymbol{\theta})$ , i.i.d. data:

$$\hat{\boldsymbol{\theta}}_{\text{MAP}} = \arg \max_{\boldsymbol{\theta}} p(\mathcal{D} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta} \in \Gamma} \log p(\boldsymbol{\theta}) + \sum_{i=1}^N \log p(\mathbf{x}^{(i)} | \boldsymbol{\theta})$$

Assuming Gaussian prior is equivalent to  $l_2$  regularization (effective  $\lambda = \sigma^2 / \alpha$  where  $\sigma^2$  is the variance of data,  $\alpha$  the variance of prior). Assuming Laplacian prior leads to  $l_1$  regularization.

## 2 Unsupervised Learning

**Dimensionality reduction:** Given a dataset  $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$  where  $\mathbf{x}^{(i)} \in \mathbb{R}^D$ , find a mapping  $f: \mathbb{R}^D \rightarrow \mathbb{R}^d$  where  $d < D$  so key characteristics of the data are retained. Can be used to save computational effort (as pre-processing), reduce overfitting, and for visualization of high-dimensional data.

### 2.1 Principal Component Analysis

Linear model  $\mathbf{z} = \mathbf{U}^T(\mathbf{x} - \mathbf{b})$  where where  $\mathbf{U} \in \mathbb{R}^{D \times d}$  is an orthonormal matrix and  $\mathbf{b} \in \mathbb{R}^D$ .  $\mathbf{z}$  is the **representation/code**; the projection of  $\mathbf{x}$  into  $\mathcal{S} = \text{col } \mathbf{U}$  is the **reconstruction**  $\hat{\mathbf{x}}$ .

$$\mathbf{b} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} \quad \Sigma = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^{(i)} - \mathbf{b})(\mathbf{x}^{(i)} - \mathbf{b})^T$$

The  $k$ -th principal component (column of  $\mathbf{U}$ ) is the eigenvector of  $\Sigma$  corresponding to the  $k$ -th largest eigenvalue. All principal components are orthogonal. Code vectors produced by PCA are de-correlated (diagonal covariance). Minimizing reconstruction error is equivalent to maximizing variance of code vectors. The first principal component is in the direction of maximum variation, then every other component is normal to all previous ones and accounts for as much of the remaining variation as possible. If  $\mathbf{X}$  is centered (subtract mean from each row), take SVD  $\mathbf{X} = \mathbf{U}_1 \mathbf{S}_1 \mathbf{V}_1^T$ , then columns of  $\mathbf{V}$  contain principal components and  $\mathbf{S}_1^2 / N$  are eigenvalues of  $\Sigma$ .

## 3 Optimization

Given  $f(\boldsymbol{\theta})$  where  $\boldsymbol{\theta} \in \mathbb{R}^n$  and  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  find  $\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$ . Define:

$$\nabla f(\boldsymbol{\theta}) = \mathbf{g}(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial f}{\partial \theta_1} \\ \vdots \\ \frac{\partial f}{\partial \theta_n} \end{bmatrix} \quad \nabla^2 f(\boldsymbol{\theta}) = \mathbf{H}(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial^2 f}{\partial \theta_1^2} & \dots & \frac{\partial^2 f}{\partial \theta_n \partial \theta_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial \theta_1 \partial \theta_n} & \dots & \frac{\partial^2 f}{\partial \theta_n^2} \end{bmatrix}$$

Karush-Kuhn-Tucker (KKT) conditions:

- First-order condition:** If  $\boldsymbol{\theta}^*$  is a minimum, then  $\nabla f(\boldsymbol{\theta}^*) = \mathbf{0}$  if it exists.
- Second-order condition:** If  $\boldsymbol{\theta}^*$  is a minimum, then  $\mathbf{H}(\boldsymbol{\theta}^*)$  is positive semi-definite, if it exists.
- Sufficient conditions:** If  $\nabla f(\boldsymbol{\theta}^*) = \mathbf{0}$  and  $\mathbf{H}(\boldsymbol{\theta}^*)$  is positive definite and both exist, then  $\boldsymbol{\theta}^*$  is a local minimum.

### 3.1 Unconstrained Optimization

Start with  $k = 0$  and an initial guess  $\boldsymbol{\theta}_0$ , repeat:

- Test for convergence of  $\boldsymbol{\theta}_k$ .
- Compute new search direction  $\mathbf{p}_k$  s.t.  $\mathbf{p}_k^T \mathbf{g}_k(\boldsymbol{\theta}_k) < 0$ .
- Compute step length  $\alpha_k > 0$  s.t.  $f(\boldsymbol{\theta}_k + \alpha_k \mathbf{p}_k) < f(\boldsymbol{\theta}_k)$ .
- Update  $\boldsymbol{\theta}_{k+1} \leftarrow \boldsymbol{\theta}_k + \alpha_k \mathbf{p}_k$ ,  $k \leftarrow k + 1$ .

**Stopping criteria:** gradient tolerance  $\epsilon_g$ , absolute tolerance  $\epsilon_a$ , relative tolerance  $\epsilon_r$ . Stop if  $\|\mathbf{g}(\boldsymbol{\theta}_k)\|_2 \leq \epsilon_g$  or  $|f(\boldsymbol{\theta}_{k+1}) - f(\boldsymbol{\theta}_k)| \leq \epsilon_a + \epsilon_r |f(\boldsymbol{\theta}_k)|$  for two successive iterations.

**Search direction:**  $\mathbf{p}_k = -\mathbf{B}_k \mathbf{g}_k$  is always valid. Steepest descent:  $\mathbf{B} = \mathbf{I}$ ; Newton's method:  $\mathbf{B} = \mathbf{H}_k^{-1}$ ; quasi-Newton:  $\mathbf{B} \approx \mathbf{H}_k^{-1}$ .

**Armijo sufficient decrease condition:**

$$f(\boldsymbol{\theta}_k + \alpha_k \mathbf{p}_k) \leq f(\boldsymbol{\theta}_k) + \mu_1 \alpha_k g_k^T \mathbf{p}_k$$

where  $\mu_1$  is typically  $\sim 10^{-4}$ . **Backtracking line search:** start with  $\alpha \in (0, 1)$ , check Armijo, and if not satisfied,  $\alpha \leftarrow \rho \alpha$  and repeat (typically  $\rho \in [0.1, 0.5]$ ).

**Steepest descent:** Use  $\mathbf{p}_k = -\frac{\mathbf{g}(\boldsymbol{\theta}_k)}{\|\mathbf{g}(\boldsymbol{\theta}_k)\|_2}$  at each step. Can show  $\mathbf{p}_k \perp \mathbf{p}_{k-1}$  so it zig zags (inefficient). Linear convergence:  $\lim_{k \rightarrow \infty} \frac{|f(\boldsymbol{\theta}_{k-1}) - f(\boldsymbol{\theta}^*)|}{|f(\boldsymbol{\theta}_k) - f(\boldsymbol{\theta}^*)|} = K$ .

**Conjugate gradient methods:** Use  $\mathbf{p}_k = -\mathbf{g}_k + \beta_k \mathbf{p}_{k-1}$  for  $k \geq 1$ , where  $\beta_k = \frac{\mathbf{g}_k^T \mathbf{g}_k}{\mathbf{g}_{k-1}^T \mathbf{g}_{k-1}}$  (Fletcher-Reeves)  $\beta_k = \frac{\mathbf{g}_k^T (\mathbf{g}_k - \mathbf{g}_{k-1})}{\mathbf{g}_{k-1}^T (\mathbf{g}_k - \mathbf{g}_{k-1})}$  (Polak-Ribieeve). More efficient due to less zig zag and still only uses first-order information.

**Newton's method:** Use  $\mathbf{p}_k = \mathbf{H}_k^{-1} \mathbf{g}_k$  at each step. Step direction only valid if  $\mathbf{H}_k$  is SPD. Quadratic convergence:  $\lim_{k \rightarrow \infty} \frac{|f(\boldsymbol{\theta}_{k-1}) - f(\boldsymbol{\theta}^*)|}{|f(\boldsymbol{\theta}_k) - f(\boldsymbol{\theta}^*)|} = K > 0$  if  $\boldsymbol{\theta}_0$  close to  $\boldsymbol{\theta}^*$  and  $\mathbf{H}(\boldsymbol{\theta}^*) > 0$ . Possible failure and numerical issues if  $\mathbf{H}_k$  is non-SPD or quadratic approximation is poor.

**Quasi-Newton methods:** Use approximate  $\mathbf{B}^{-1} \approx \mathbf{H}_k^{-1}$  and update approximate Hessian as  $\mathbf{B}_{k+1}^{-1} = \mathbf{B}_k^{-1} + \Delta \mathbf{B}_k$  at each step. Convergence between linear and quadratic.

**SR1 update:** Update  $\mathbf{B}_k$  by a symmetric rank-1 matrix at each step.

$$\mathbf{B}_{k+1}^{-1} = \mathbf{B}_k^{-1} + \frac{(\mathbf{s}_k - \mathbf{B}_k^{-1} \mathbf{y}_k)(\mathbf{s}_k - \mathbf{B}_k^{-1} \mathbf{y}_k)^T}{(\mathbf{s}_k - \mathbf{B}_k^{-1} \mathbf{y}_k)^T \mathbf{y}_k}$$

### 3.2 Unconstrained Optimization

Minimize  $f(\boldsymbol{\theta})$  subject to constraints  $g_i(\boldsymbol{\theta}) \geq 0, h_j(\boldsymbol{\theta}) = 0$  where  $i = 1, 2, \dots, m$  and  $j = 1, 2, \dots, q$  and  $\boldsymbol{\theta}_i \leq \boldsymbol{\theta} \leq \boldsymbol{\theta}_m$ .

**Penalty methods:** Minimize  $\pi(\boldsymbol{\theta}, \rho_k) = f(\boldsymbol{\theta}) + \rho_k \phi(\boldsymbol{\theta})$ , where  $\phi(\boldsymbol{\theta})$  is ero when constraints are met, positive when constraints violated, e.g. quadratic penalty function  $\phi(\boldsymbol{\theta}) = \sum_{i=1}^m (\max(0, -g_i(\boldsymbol{\theta})))^2 + \sum_{j=1}^q (h_j(\boldsymbol{\theta}))^2$ . At each step, check convergence conditions, find  $\boldsymbol{\theta}_{k+1} = \arg \min_{\boldsymbol{\theta} \in \Gamma} \pi(\boldsymbol{\theta}, \rho_k)$ , increment

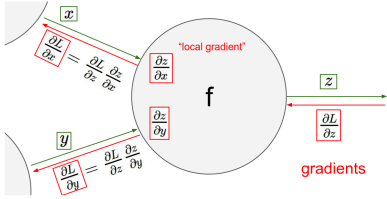
**Logistic sigmoid:** for binary classification:  $\sigma(z) = \frac{1}{1+e^{-z}}$   
**Softmax:** for multi-class classification:  $\text{softmax}(\mathbf{z}) = \frac{e^{z_j}}{\sum_k e^{z_k}}$   
 $\text{logsumexp}(\mathbf{z}) = \log(\sum_i e^{z_i}) \implies \log \text{softmax}(\mathbf{z})_j = z_j - \text{logsumexp}(\mathbf{z})$   
**Classification dataset likelihoods:**  
 Binary: i.i.d. Bernoulli,  $\Pr(y|\mathbf{x}, \mathbf{w}) = \hat{f}(\mathbf{x}; \mathbf{w})^{y(1-\hat{f}(\mathbf{x}; \mathbf{w}))^{1-y}}$

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} \sum_{i=1}^N y_i^{(i)} \log f(\mathbf{x}^{(i)}; \mathbf{w}) + (1 - y_i^{(i)}) \log(1 - f(\mathbf{x}^{(i)}; \mathbf{w}))$$

Multiclass: i.i.d. Categorical,  $\Pr(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \prod_{j=1}^K \hat{f}_j(\mathbf{x}; \mathbf{w})^{y_j}$

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} \sum_{i=1}^N \sum_{j=1}^K y_j^{(i)} \log f_j(\mathbf{x}^{(i)}; \mathbf{w})$$

**Backpropagation:** Each node  $z = f(x_1, \dots, x_N)$  in the computational graph is passed an upstream gradient  $\frac{\partial L}{\partial z}$  (sum of all downstream gradients of the nodes that  $z$  connects to) and computes local gradients  $\frac{\partial z}{\partial x_i}$ ; pass gradients downstream as  $\frac{\partial L}{\partial x_i} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x_i}$ ; full gradient is obtained when leaf is reached. +: distribute gradients; max/min: route gradient to 1 input;  $\times$ : swap gradients.



**Weight initialization:** Aim for zero mean, unit variance of initial hidden states to avoid vanishing or exploding gradients. *Xavier:* sample from unit Gaussian, divide by square root of number of neuron inputs; *ReLU:* divide by square root of half the number of inputs.

#### 4.1 Regularization Techniques

**Weight decay:** Equivalent to  $l_2$  regularization; keeps weights small.  
**Early stopping:** Stopping when validation loss increases, before convergence in (training) gradient.  
**Data augmentation:** Creating new samples from existing ones through class-preserving transforms, e.g. rotation, blurring, cropping, etc. or adversarial training samples designed to mislead the model.  
**Bagging/bootstrap aggregation:** Generate  $k$  datasets of equivalent size from  $\mathcal{D}$  by sampling with replacement, and train  $k$  models on datasets; at inference time, average predictions (regression) or vote on class (classification).  
**Dropout:** During training, drop weights randomly with probability  $1 - \pi$  (typ.  $\pi \in [0.5, 0.8]$ ) at each SGD iteration; during inference, scale weights by  $\pi$ . Equivalent to bagging with a  $k$  that increases exponentially with network size (but same model across all datasets).  
**Weight sharing:** Using prior knowledge to constrain weights to be similar; CNNs are an example.

#### 4.2 Network Architectures

**Convolutional Neural Networks (CNNs):** Convolving 2D (or higher) filters with the input to produce output for the next layer, keeping structure. Use on highly structured data, e.g. images. Pooling layers between convolutional layers perform spatial downsampling (but keeps depth). Layers will learn progressively higher-level features until linearly separable data at the last layer (or passed to MLP).  
**Autoencoders:** Unsupervised model consisting of an encoder from high-dimensional input space to low-dimensional feature space, and a decoder from feature space back to output space. During training the model is made to reconstruct its input. Can be used for dimensionality reduction, new sample generation, noise filtering, anomaly detection, etc. or trained further on labelled data for final task (semi-supervised learning).

### 5 Bayesian Inference

Given the **likelihood**  $p(y^{(1)}, \dots, y^{(N)}|\theta) = p(\mathcal{D}|\theta)$  and the **prior distribution**  $p(\theta)$ , we wish to find:

- Posterior distribution:**  $p(\theta|\mathcal{D}) = \frac{p(\mathcal{D}|\theta)p(\theta)}{\int p(\mathcal{D}|\theta')p(\theta')d\theta'}$
- Posterior predictive distribution:**  $p(y'|\mathcal{D}) = \int p(y'|\theta)p(\theta|\mathcal{D})d\theta$

**Conjugate prior:** A prior with similar mathematical form as the likelihood; e.g. beta for Bernoulli/geometric, Gaussian for Gaussian, Dirichlet for multinomial, gamma for Poisson/exponential.  
**Evidence/marginal likelihood:** Likelihood of observing the data when marginalizing across all possible parameters  $\mathbf{w}$  according to the prior:

$$p(\mathbf{y}|\mathbf{X}) = \int p(\mathbf{y}|\mathbf{w}, \mathbf{X})p(\mathbf{w})d\mathbf{w}$$

$$p(\mathbf{y}|\mathbf{X}, \alpha, \sigma^2) = \int p(\mathbf{y}|\mathbf{w}, \mathbf{X}, \sigma^2)p(\mathbf{w}|\alpha)d\mathbf{w}$$

#### 5.1 Bayesian Linear Regression

Given dataset  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ , assume  $y \sim \mathcal{N}(y|\mathbf{w}^T\phi(\mathbf{x}), \sigma^2)$  giving likelihood  $\log p(\mathbf{y}|\mathbf{w}, \mathbf{X}, \sigma^2) = \sum_{i=1}^N \log \mathcal{N}(y^{(i)}|\mathbf{w}^T\phi(\mathbf{x}^{(i)}), \sigma^2)$ , with prior  $p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha\mathbf{1})$ . Bayesian linear regression gives:

- **Weight Posterior:**  $p(\mathbf{w}|\mathcal{D}) = \mathcal{N}(\mathbf{w}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- **Posterior predictive:**  $p(y'|\mathbf{x}', \mathcal{D}) = \mathcal{N}(y'|\boldsymbol{\mu}^T\phi(\mathbf{x}'), \phi(\mathbf{x}')^T\boldsymbol{\Sigma}\phi(\mathbf{x}') + \sigma^2)$  where

$$\boldsymbol{\mu} = \frac{1}{\sigma^2}\boldsymbol{\Sigma}\boldsymbol{\Phi}^T\mathbf{y} \quad \boldsymbol{\Sigma} = \left(\frac{1}{\sigma^2}\boldsymbol{\Phi}^T\boldsymbol{\Phi} + \frac{1}{\alpha}\mathbf{1}\right)^{-1}$$

**Complexity:**  $\mathcal{O}(NM^2 + M^3)$  time,  $\mathcal{O}(NM + M^2)$  memory.  
**Type-II inference:** Estimate  $\sigma^2$  and  $\alpha$  through numerical maximization of the (log) evidence/marginal likelihood (where  $N$  is number of samples):

$$\log p(\mathbf{y}|\mathbf{X}, \alpha, \sigma^2) = -\frac{N}{2}\log\alpha - \frac{N}{2}\log(\sigma^2) - \frac{1}{2\sigma^2}\mathbf{y}^T\mathbf{y}$$

$$+ \frac{1}{2}\boldsymbol{\mu}^T\boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + \frac{1}{2}\log\det(\boldsymbol{\Sigma}) - \frac{N}{2}\log(2\pi)$$

#### 5.2 Gaussian Processes

Define: **Gram matrix**  $\mathbf{K}_{\mathcal{X}, \mathcal{X}} = \alpha\boldsymbol{\Phi}\boldsymbol{\Phi}^T \in \mathbb{R}^{N \times N}$ , where entry  $ij$  is  $[\mathbf{K}_{\mathcal{X}, \mathcal{X}}]_{ij} = \alpha\phi^T(\mathbf{x}^{(i)})\phi(\mathbf{x}^{(j)}) = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)})$ ;  $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is the **kernel**. Let  $\mathbf{k}_{\mathcal{X}, \mathbf{x}'} = [k(\mathbf{x}^{(1)}, \mathbf{x}') \dots k(\mathbf{x}^{(N)}, \mathbf{x}')]^T$ . Then the kernelized posterior predictive is

$$p(y'|\mathbf{x}', \mathcal{D}) = \mathcal{N}(y'|\mu_p, \sigma_p^2) \quad \mu_p = \mathbf{k}_{\mathbf{x}', \mathcal{X}}(\mathbf{K}_{\mathcal{X}, \mathcal{X}} + \sigma^2\mathbf{1})^{-1}\mathbf{y}$$

$$\sigma_p^2 = k(\mathbf{x}', \mathbf{x}') - \mathbf{k}_{\mathbf{x}', \mathcal{X}}(\mathbf{K}_{\mathcal{X}, \mathcal{X}} + \sigma^2\mathbf{1})^{-1}\mathbf{k}_{\mathcal{X}, \mathbf{x}'} + \sigma^2$$

$$\text{cov}(y_1', y_2') = k(\mathbf{x}_1', \mathbf{x}_2') - \mathbf{k}_{\mathbf{x}_1', \mathcal{X}}(\mathbf{K}_{\mathcal{X}, \mathcal{X}} + \sigma^2\mathbf{1})^{-1}\mathbf{k}_{\mathcal{X}, \mathbf{x}_2'} + \sigma^2$$

which is equivalent to a GLM with  $l_2$  error and regularization,  $\lambda = \sigma^2/\alpha$ .  $k(\cdot, \cdot)$  specifies covariance between two observations; any collection of observations are jointly Gaussian, making this a **Gaussian process**.

**Complexity:**  $\mathcal{O}(N^3)$  time,  $\mathcal{O}(N^2)$  memory. Much more efficient than normal Bayesian linear regression when  $M \gg N$  or possibly infinite (same kernels as GLMs).

**Combining kernels:** Through summation  $k(x, y) = k_1(x, y) + k_2(x, y)$  or  $k(\mathbf{x}, \mathbf{y}) = k_1(x_1, y_1) + k_2(x_2, y_2)$ , multiplication  $k(x, y) = k_1(x, y)k_2(x, y)$  or  $k(\mathbf{x}, \mathbf{y}) = k_1(x_1, y_1)k_2(x_2, y_2)$ , or composition  $k(x, y) = k_1(f(x), f(y))$ , all preserving positive-definiteness.

**Type-II inference:** In function-space,

$$\log p(\mathbf{y}|\mathbf{X}, \alpha, \sigma^2) = -\frac{N}{2}\log(2\pi) - \frac{1}{2}\log\det(\mathbf{K}_{\mathcal{X}, \mathcal{X}} + \sigma^2\mathbf{1})$$

$$- \frac{1}{2}\mathbf{y}^T(\mathbf{K}_{\mathcal{X}, \mathcal{X}} + \sigma^2\mathbf{1})^{-1}\mathbf{y}$$

#### 5.3 Approximate Bayesian Inference

Given observed evidence,  $X_E$  and unobserved variables,  $X_F$ , find

$$p(X_F|X_E) = \frac{p(X_E, X_F)}{p(X_E)}$$

$p(X_E, X_F) = p(X_E|X_F)p(X_F)$  is often known from the model and prior assumptions, but  $p(X_E) = \int p(X_E, X_F)dX_F$  is intractable to compute.

**Quadrature:** Numerically integrate  $p(X_E) = \int p(X_E, X_F)dX_F$  by sampling a finite number of points. Infeasible due to the exponential growth in points required. Error usually scales as  $\mathcal{O}(n^{-\frac{1}{d}})$  where  $d$  is the dimensionality of the integral. Beyond  $d \geq 3$  Monte Carlo is better.

**Laplace approximation:** Let  $X_F = \mathbf{z}$  and  $\hat{\mathbf{z}}_{\text{MAP}} = \arg \max_{\mathbf{z}} \log \tilde{p}(\mathbf{z})$  where  $\tilde{p}(\mathbf{z}) = p(X_E, \mathbf{z})$ , then  $p(\mathbf{z}|X_E) = \frac{1}{2}p(X_E, \mathbf{z}) = \frac{1}{2}\tilde{p}(\mathbf{z})$ . Then:

$$p(\mathbf{z}|X_E) \approx \mathcal{N}(\mathbf{z}|\hat{\mathbf{z}}_{\text{MAP}}, \mathbf{A}^{-1}) \quad \mathbf{A} = -\nabla^2 \log \tilde{p}(\mathbf{z})|_{\mathbf{z}=\hat{\mathbf{z}}_{\text{MAP}}}$$

Note: from a second-order Taylor expansion:

$$\log p(\mathbf{z}|X_E) \approx \log \tilde{p}(\hat{\mathbf{z}}_{\text{MAP}}) - \frac{1}{2}(\mathbf{z} - \hat{\mathbf{z}}_{\text{MAP}})^T \mathbf{A}(\mathbf{z} - \hat{\mathbf{z}}_{\text{MAP}}) + \text{const.}$$

Simple but poor performance, since it does not account for global information. However, if posterior is Gaussian, this is exact. Requires existence of Hessian for likelihood and prior, so certain distributions (e.g. Laplace) might be bad.

#### 5.3.1 Stochastic Variational Inference (SVI)

**Kullback-Leibler (KL) divergence:** For two distributions  $p(\mathbf{z})$  and  $q(\mathbf{z})$ ,

$$KL(q\|p) = \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[ \log \frac{q(\mathbf{z})}{p(\mathbf{z})} \right] \implies \begin{cases} KL(q\|p) \geq 0 \\ KL(q\|p) = 0 \iff q = p \\ KL(q\|p) \neq KL(p\|q) \end{cases}$$

**Reverse-KL (information projection):** take  $KL(q\|p)$  which penalizes  $q$  having mass where  $p$  has none (compressing  $q$  to fit to one peak of  $p$ ); **forward-KL (moment projection):** take  $KL(p\|q)$  which penalizes  $q$  missing mass where  $p$  has mass (stretching out  $q$  to cover all peaks of  $p$ ). Reverse KL commonly used for computational reasons.

**Stochastic Variational Inference (SVI):** Approximating the true conditional  $p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{z}, \mathbf{x})}{p(\mathbf{x})}$  by a simpler distribution,  $q(\mathbf{z}|\theta)$ , from a known family of distributions, such that  $KL(q(\mathbf{z}|\theta)\|p(\mathbf{z}|\mathbf{x}))$  is minimized.

**Evidence lower bound (ELBO):**  $\mathcal{L}(\theta, \mathbf{x}) = -\mathbb{E}_{\mathbf{z} \sim q} \left[ \log \frac{q(\mathbf{z}|\theta)}{p(\mathbf{z}, \mathbf{x})} \right]$  is a lower bound for  $\log p(\mathbf{x})$ ; maximizing ELBO minimizes KL. Has gradient

$$\nabla_{\theta} \mathcal{L}(\theta, \mathbf{x}) = \nabla_{\theta} \int q(\mathbf{z}|\theta) \log \frac{p(\mathbf{z}, \mathbf{x})}{q(\mathbf{z}|\theta)} d\mathbf{z}$$

**Score function/REINFORCE gradient estimator:** By Monte Carlo:

$$\nabla_{\theta} \mathcal{L}(\theta, \mathbf{x}) = \mathbb{E}_{\mathbf{z} \sim q} \left[ \nabla_{\theta} \log q(\mathbf{z}|\theta) \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\theta)} \right]$$

$$\approx \frac{1}{B} \sum_{i=1}^B \nabla_{\theta} \log q(\mathbf{z}^{(i)}|\theta) \log \frac{p(\mathbf{x}, \mathbf{z}^{(i)})}{q(\mathbf{z}^{(i)}|\theta)}$$

**Pathwise/reparametrization gradient estimator:** Find  $T(\boldsymbol{\varepsilon}, \theta)$  such that for  $\boldsymbol{\varepsilon} \sim p(\boldsymbol{\varepsilon})$ , then  $\mathbf{z} = T(\boldsymbol{\varepsilon}, \theta) \implies \mathbf{z} \sim q(\mathbf{z}|\theta)$ . Then Monte Carlo:

$$\nabla_{\theta} \mathcal{L}(\theta, \mathbf{x}) = \mathbb{E}_{\boldsymbol{\varepsilon} \sim p(\boldsymbol{\varepsilon})} \left[ \nabla_{\theta} \log \frac{p(\mathbf{x}, T(\boldsymbol{\varepsilon}, \theta))}{q(T(\boldsymbol{\varepsilon}, \theta)|\theta)} \right]$$

$$\approx \frac{1}{B} \sum_{i=1}^B \nabla_{\theta} \log \frac{p(\mathbf{x}, T(\boldsymbol{\varepsilon}^{(i)}, \theta))}{q(T(\boldsymbol{\varepsilon}^{(i)}, \theta)|\theta)}$$

e.g. for a Gaussian,  $\theta = \{\mu, \sigma^2\}$ , let  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\boldsymbol{\varepsilon}|0, 1)$  and  $T(\boldsymbol{\varepsilon}, \theta) = \sigma\boldsymbol{\varepsilon} + \mu$ , then  $\mathbf{z} \sim \mathcal{N}(\mathbf{z}|\mu, \sigma^2)$ .

#### 5.3.2 Expectation Approximation

Instead of  $p(\mathbf{x})$  sometimes we only need  $I = \mathbb{E}_{\mathbf{x} \sim p(\mathbf{x})} [\phi(\mathbf{x})]$ .

**Monte Carlo:** Collect  $R$  random samples from  $p(\mathbf{x})$ , denoted  $\mathbf{x}^{(i)}$ , then

$$I = \int \phi(\mathbf{x})p(\mathbf{x})d\mathbf{x} \approx \hat{I} = \frac{1}{R} \sum_{i=1}^R \phi(\mathbf{x}^{(i)})$$

Has standard deviation proportional to  $1/\sqrt{R}$ , independent of dimension; unbiased estimator. For dimension  $d \geq 3$  this generally outperforms quadrature. However, we might only know the unnormalized  $\tilde{p}(\mathbf{x}) = Zp(\mathbf{x})$ , or sampling may be difficult due to high dimensionality, in which case we need importance sampling.

**Importance Sampling:** Let  $q(\mathbf{x})$  be the *sampler density*, a density easy to sample from; then for some unnormalized  $\tilde{p}(\mathbf{x}) = Zp(\mathbf{x})$ ,

$$I = \frac{\mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[ \frac{\phi(\mathbf{x})\tilde{p}(\mathbf{x})}{q(\mathbf{x})} \right]}{\mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[ \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})} \right]} \implies \hat{I} = \frac{\frac{1}{R} \sum_{r=1}^R \frac{\phi(\mathbf{x}^{(r)})\tilde{p}(\mathbf{x}^{(r)})}{q(\mathbf{x}^{(r)})}}{\frac{1}{R} \sum_{r=1}^R \frac{\tilde{p}(\mathbf{x}^{(r)})}{q(\mathbf{x}^{(r)})}} = \frac{\sum_r w_r \phi(\mathbf{x}^{(r)})}{\sum_r w_r}$$

where each  $w_r = \tilde{p}(\mathbf{x}^{(r)})/q(\mathbf{x}^{(r)})$  is an *importance weight*, which is larger if sampling from  $q(\mathbf{x})$  tends to under-represent the point, or smaller for over-represent. Note if  $\tilde{p}(\mathbf{x}) = p(\mathbf{x})$  is normalized, then

$$\mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[ \frac{\tilde{p}(\mathbf{x})}{q(\mathbf{x})} \right] = 1 \implies \hat{I} = \frac{1}{R} \sum_{r=1}^R \frac{\phi(\mathbf{x}^{(r)})\tilde{p}(\mathbf{x}^{(r)})}{q(\mathbf{x}^{(r)})} = \frac{1}{R} \sum_r w_r \phi(\mathbf{x}^{(r)})$$

Sampler should have heavy tails (e.g. Cauchy instead of Gaussian) to compensate for differences between  $p(\mathbf{x})$  and  $q(\mathbf{x})$ . If the sampler is chosen improperly, the variance of the result can be extremely high. In high dimensions, if the sampler distribution is not a near-perfect approximation of the target, then the entire sum will likely be dominated by a few samples with a huge weight, leading to a very bad estimate.

### 6 Useful Identities

**Matrix differentiation:**  $\frac{\partial}{\partial \mathbf{z}} (\mathbf{z}^T \mathbf{A} \mathbf{z}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{z} \quad \frac{\partial}{\partial \mathbf{z}} (\mathbf{A} \mathbf{z}) = \mathbf{A}^T$

**Laplace distribution:**  $\text{Lap}(c|\mu, b) = \frac{1}{2b} e^{-\frac{|c-\mu|}{b}}$  with  $\sigma^2 = 2b^2$ .

**Gaussian distribution:**  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$ ; multi-var ver.:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{D}{2}} \det(\boldsymbol{\Sigma})^{\frac{1}{2}}} \exp \left[ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right]$$

Partition Gaussian  $\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  as:

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}, \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \implies \mathbf{x} \sim f(\mathbf{x}_1, \mathbf{x}_2)$$

Then  $\mathbf{x}_1 \sim \mathcal{N}(\mathbf{x}_1|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11})$ ,  $\mathbf{x}_2 \sim \mathcal{N}(\mathbf{x}_2|\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})$ , and

$$f(\mathbf{x}_1|\mathbf{x}_2) = \mathcal{N}(\mathbf{x}_1|\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21})$$

**Beta distribution:**  $p(\theta|\alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1}(1-\theta)^{\beta-1}$  ( $\alpha = \beta = 1$  gives uniform dist.) with mean  $\frac{\alpha}{\alpha+\beta}$  and mode  $\frac{\alpha-1}{\alpha+\beta-2}$ .

**Dirichlet distribution:**  $p(\theta|\alpha_1, \dots, \alpha_K) = \frac{\Gamma(\sum_k \alpha_k)}{\prod_k \Gamma(\alpha_k)} \prod_k \theta_k^{\alpha_k-1}$  where  $\alpha_k > 0$ ,  $\sum_k \alpha_k = 1$ .

**Sigmoid & Derivative:**  $\sigma(x) = \frac{1}{1+e^{-x}} \implies \frac{d}{dx} \sigma(x) = \sigma(x)(1-\sigma(x))$