

Background

Linear Algebra

$$\begin{aligned}\|\mathbf{x}\|_p &= (\sum |x_i|^p)^{1/p} & \|\mathbf{x}\|_\infty &= \max |x_i| \\ \text{tr}(\mathbf{A}\mathbf{x}\mathbf{x}^T) &= \mathbf{x}^T \mathbf{A} \mathbf{x} \\ |\mathbf{A}\mathbf{B}| &= |\mathbf{A}||\mathbf{B}| & |\mathbf{A}^m| &= |\mathbf{A}|^m \\ (\mathbf{A} + \mathbf{UCV})^- &= \mathbf{A}^- - \mathbf{A}^- \mathbf{U}(\mathbf{C}^+ + \mathbf{V}\mathbf{A}^+ \mathbf{U})^- \mathbf{V}\mathbf{A}^- \\ (\mathbf{A} + \mathbf{B})^{-1} &= \mathbf{A}^{-1} - \mathbf{A}^{-1}(\mathbf{A} + \mathbf{B})^{-1}\mathbf{A}^{-1} \\ \mathbf{U}(\mathbf{V}\mathbf{U} + \mathbf{I})^{-1} &= (\mathbf{U}\mathbf{V} + \mathbf{I})^{-1}\mathbf{U} \\ \mathbf{I} - \mathbf{A}(\mathbf{I} + \mathbf{A})^{-1} &= (\mathbf{I} + \mathbf{A})^{-1}\end{aligned}$$

Probability

$$\text{Ber}(x|\theta) = \theta^x(1-\theta)^{1-x} \quad 0 \leq \theta \leq 1$$

$$\mathbb{P}[X|Y] = \frac{\mathbb{P}[X,Y]}{\mathbb{P}[Y]} = \frac{\mathbb{P}[Y|X]\mathbb{P}[X]}{\mathbb{P}[Y]}$$

$$p_Y(y) = p_X(g^{-1}(y)) \left| \det \frac{\partial g^{-1}(y)}{\partial y} \right|$$

$$\begin{aligned}\mathbb{E}[X] &= \int_{\Omega} xp(x) dx = \int_{\omega} x\mathbb{P}[X=x] dx \\ \mathbb{E}_{Y|X}[Y] &= \mathbb{E}_Y[Y|X] \mathbb{I}_{\mathcal{E}_Y}[\mathbb{E}_X[X|Y]] = \mathbb{E}_X[X] \\ \mathbb{E}_{X,Y}[f(X,Y)] &= \mathbb{E}_X \mathbb{E}_{Y|X}[f(X,Y)|X]\end{aligned}$$

$$\begin{aligned}\mathbb{V}(X) &= \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 \\ \text{Var}(X) &= \mathbb{E}[\text{Var}(X|Y)] + \text{Var}(\mathbb{E}[X|Y]) \\ \mathbb{V}[X+Y] &= \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}(X,Y) \\ \text{Cov}(\mathbf{X}, \mathbf{Y}) &= \mathbb{E}[\mathbf{X}\mathbf{Y}^T] - \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{Y}]^T \\ \text{Cov}(\mathbf{A}\mathbf{X} + c, \mathbf{B}\mathbf{Y} + d) &= \mathbf{A}\text{Cov}(\mathbf{X}, \mathbf{Y})\mathbf{B}^T\end{aligned}$$

$$\mathcal{N}(x|\mu, \Sigma) = \frac{\exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu))}{(2\pi)^{D/2} |\Sigma|^{1/2}}$$

$$X = \Sigma^{1/2} \mathcal{N}(0, 1) + \mu \sim \mathcal{N}(\mu, \Sigma)$$

$$Y = MX + b \sim \mathcal{N}(M\mu + b, M\Sigma M^T)$$

$$X, Y \stackrel{\text{iid}}{\sim} \mathcal{N}: \quad X+Y \sim \mathcal{N}(\mu+\mu', \Sigma+\Sigma')$$

Conditional Gaussian

$$P\left(\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix}; \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}\right)$$

$$p(\mathbf{Y}|\mathbf{X} = \mathbf{x}) = \mathcal{N}(\mu, \Sigma)$$

$$\mu = \mu_2 + \Sigma_{21}\Sigma_{11}^{-1}(\mathbf{x} - \mu_1)$$

$$\Sigma = \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}$$

Inequalities and Estimators

$$\text{Jensen:} \quad \log(\sum_i \lambda_i^{(\geq 0)} x_i) \geq \sum_i \lambda_i \log(x_i)$$

$$\text{Chebyshev: } \mathbb{P}(|\hat{X} - X| \geq \varepsilon) \leq \frac{MSE[\hat{X}]}{\varepsilon^2}$$

$$\text{Estimators:} \quad \text{Unbiased: } \mathbb{E}[\hat{\theta}] = \theta^*$$

$$\text{Consistent: } \mathbb{P}(|\hat{\theta} - \theta^*| < \varepsilon) \rightarrow 0 \text{ convP}$$

$$\text{Asymp Normal: } (\hat{\theta} - \theta^*)\hat{se}^{-1} \sim \mathcal{N}(0, 1)$$

$$\text{Rao-Cra.: } \mathbb{E}_{x|\theta}[(\theta - \hat{\theta})^2] \geq \frac{(\frac{\partial}{\partial \theta} b_{\theta+1})^2}{\mathbb{E}_{x|\theta}[\Lambda^2]} + b_{\theta}^2$$

$$b_{\hat{\theta}} = \mathbb{E}_{x|\theta}[\hat{\theta}] - \theta \quad \Lambda = \frac{\partial}{\partial \theta} \log p(x|\theta)$$

$$\mathbb{E}_{x|\theta}[\Lambda] = 0 \rightarrow \mathbb{E}_{x|\theta}[\Lambda \hat{\theta}] = \frac{\partial}{\partial \theta} b_{\theta+1} + 1$$

$$\rightarrow \text{Cov}(\Lambda, \hat{\theta}) \rightarrow \text{Cauchy}$$

$$\text{Var}[\hat{\theta}] \geq \mathcal{I}_n(\theta)^{-1} = -\mathbb{E}\left[\frac{\partial^2 \log p[\mathcal{X}_n|\theta]}{\partial \theta^2}\right]^{-1}$$

$$\text{Efficiency of } \hat{\theta}: e(\theta_n) = \frac{1}{\text{Var}[\hat{\theta}_n|\mathcal{I}_n(\theta)]}$$

$$\hat{\theta}_{JS} = \left(1 - \frac{(d-2)\sigma^2}{\|\mathbf{y}\|^2}\right) y$$

$$\text{Derivatives: } \frac{\partial}{\partial \mathbf{x}}(\mathbf{b}^T \mathbf{x}) = \frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^T \mathbf{b}) = \mathbf{b}$$

$$\frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^T \mathbf{A} \mathbf{x}) = (\mathbf{A}^T + \mathbf{A})\mathbf{x}$$

$$\frac{\partial}{\partial \mathbf{X}}(\mathbf{c}^T \mathbf{X} \mathbf{b}) = \mathbf{c} \mathbf{b}^T$$

$$\frac{\partial}{\partial \mathbf{X}}(\|\mathbf{X}\|_F^2) = 2\mathbf{X}$$

$$\frac{\partial}{\partial \mathbf{x}} \|\mathbf{x}\|_2 = \frac{\mathbf{x}}{\|\mathbf{x}\|_2} \quad \frac{\partial}{\partial \mathbf{x}} \|f(\mathbf{x})\|_1 = \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}}^T \text{sgn}(\mathbf{x})$$

$$\frac{\partial}{\partial \mathbf{x}} (\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2) = 2\mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{b})$$

$$\frac{\partial}{\partial \mathbf{X}} (\|\mathbf{X}\|) = |\mathbf{X}| \cdot \mathbf{X}^{-1}, \quad |\mathbf{X}|^{-1} = |\mathbf{X}^{-1}|$$

$$\frac{\partial}{\partial \mathbf{X}} f(\mathbf{X})^T = \frac{\partial f(\mathbf{X})}{\partial \mathbf{X}}^T \quad \frac{\partial}{\partial \mathbf{X}} \text{tr} f(\mathbf{X}) = \text{tr} \frac{\partial f(\mathbf{X})}{\partial \mathbf{X}}$$

$$\frac{\partial}{\partial \mathbf{X}} \det f(\mathbf{X}) = \det f(\mathbf{X}) \text{tr}(f(\mathbf{X})^{-1} \frac{\partial f(\mathbf{X})}{\partial \mathbf{X}})$$

$$\frac{\partial}{\partial \mathbf{X}} f(\mathbf{X})^{-1} = -f(\mathbf{X})^{-1} \frac{\partial f(\mathbf{X})}{\partial \mathbf{X}} f(\mathbf{X})^{-1}$$

Quadratic Forms

$$\mathbf{x}^T \mathbf{A} \mathbf{x} + 2\mathbf{b}^T \mathbf{x} + c = (\mathbf{x} + \mathbf{A}^{-1}\mathbf{b})^T \mathbf{A}(\mathbf{x} + \mathbf{A}^{-1}\mathbf{b}) - \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b} + c,$$

$$ax^2 + bx + c = (x + \frac{b}{2a})^2 - (\frac{b}{2a})^2 + c$$

Information Theory

$$H[p] = \mathbb{E}_{\mathbf{x} \sim p}[-\log p(\mathbf{x})]$$

$$H[p||q] = \mathbb{E}_{\mathbf{x} \sim p}[-\log q(\mathbf{x})]$$

$$\text{KL}[p||q] = H[p||q] - H[p]$$

$$\text{KL}[p||q] = \mathbb{E}_{\theta \sim p} \left[\log \left(\frac{p(\theta)}{q(\theta)} \right) \right]$$

$$\text{KL}[p||q] \neq \text{KL}[q||p] \geq 0$$

$$H[\mathbf{X}] = \mathbb{E}_{\mathbf{X} \sim p}[-\log p(\mathbf{X})]$$

$$H[\mathbf{X}|\mathbf{Y} = y] = \mathbb{E}_{\mathbf{X} \sim p(\cdot|y)}[-\log p(\mathbf{X}|y)]$$

$$H[\mathbf{X}|\mathbf{Y}] = \mathbb{E}_y[H[\mathbf{X}|\mathbf{Y} = y]]$$

$$H[\mathbf{X}|\mathbf{Y}] = H[\mathbf{Y}|\mathbf{X}] + H[\mathbf{X}] - H[\mathbf{Y}]$$

$$H[\mathbf{X}, \mathbf{Y}] = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim p(\cdot, \cdot)}[-\log p(\mathbf{X}, \mathbf{Y})]$$

$$I[\mathbf{X}; \mathbf{Y}] = H[\mathbf{X}] - H[\mathbf{X}|\mathbf{Y}] \geq 0$$

$$I[\mathbf{X}; \mathbf{Y}|\mathbf{Z}] = I[\mathbf{X}; \mathbf{Y}, \mathbf{Z}] - I[\mathbf{X}; \mathbf{Z}]$$

$$H(\mathcal{N}(\mu, \Sigma) = \frac{1}{2} \ln(\det(2\pi e \Sigma))$$

$$\text{KL}(\mathcal{N}(a, A) || \mathcal{N}(b, B)) = \frac{1}{2} (\text{tr}(B^{-1}A) + (a-b)^T B^{-1}(a-b) - d + \ln(\frac{\det B}{\det A}))$$

Risks

$$\text{Expected Risk: } R(f) = P(f(X) \neq y)$$

$$\mathcal{R}(f) = \sum_{y \leq k} P(y) \mathbb{E}_{P(x|y)}[1_{f(x) \neq y} | Y = y]$$

$$\text{Empirical Risk Minimizer (ERM)} \hat{f}: \hat{f} \in \arg \min_{f \in \mathcal{H}} \hat{R}(\hat{f}, \mathcal{D}^{train})$$

$$\hat{R}(\hat{f}, \mathcal{D}^{train}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(Y_i, \hat{f}(X_i))$$

$$\hat{R}(\hat{f}, \mathcal{D}^{test}) = \frac{1}{m} \sum_{i=n+1}^{n+m} \mathcal{L}(Y_i, \hat{f}(X_i))$$

$$\text{Loss Fcts: } \mathcal{L}(y, z) \quad z = w^T x$$

$$\mathcal{L}^{0/i} = \mathbb{I}[\text{sign}(z) \neq y]$$

$$\mathcal{L}^{\text{hinge}} = \max(0, 1 - yz) \quad \text{for SVM's}$$

$$\mathcal{L}^{\text{percep}} = \max(0, -yz)$$

$$\mathcal{L}^{\text{logistic}} = \log(1 + \exp(-yz))$$

$$\mathcal{L}^{\text{exp}} = \exp(-yz) \quad \text{for AdaBoost}$$

$$\mathcal{L}^{\text{CE}} = -[y' \log z' + (1 - y') \log(1 - z')]$$

$$y' = \frac{1+y}{2}, \quad z' = \frac{1+z}{2}$$

Optimization

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla_{\theta} \mathcal{L} + \mu(\theta^{(t)} - \theta^{(t-1)})$$

$$GD: \quad \theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla_{\theta} \mathcal{L}$$

$$SGD: \quad \theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla \mathcal{L}(\theta^{(t)}, x_i, y_i)$$

$$NGD: \quad \theta^{(t+1)} \leftarrow \theta^{(t)} - \eta (\nabla_{\theta}^2 \mathcal{L})^{-1} \nabla_{\theta} \mathcal{L}$$

$$\rightarrow f(x+t) \approx f(x) + t f'(x) + \frac{1}{2} f''(x) t^2 = 0$$

Parametric Density Estimation

$$\text{Assume prior } \mathbb{P}(\theta),$$

$$\text{Likelihood: } \mathbb{P}[\mathcal{X}|\theta] = \prod_{i \leq n} p(x_i|\theta)$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \mathbb{P}[\mathcal{X}|\theta]$$

$$\hat{\theta}_{MAP} = \arg \max_{\theta} [P(\theta|\mathcal{X}) = P(\mathcal{X}|\theta)P(\theta)]$$

$$\text{Solve } \nabla_{\theta} \log P(\mathcal{X}|\theta)P(\theta) = 0$$

1-D Gaussian Bayesian learning

$$X|\theta \sim \mathcal{N}(\theta, \sigma^2) \quad \theta \sim \mathcal{N}(m_0, s_0^2)$$

$$\theta|X \sim \mathcal{N}(\mu_n, \sigma_n^2)$$

$$\sigma_n^2 = \frac{\sigma^2 s_0^2}{ns_0^2 + \sigma^2}, \quad \mu_n = \frac{ns_0^2 \bar{x} + m_0 \sigma^2}{ns_0^2 + \sigma^2}$$

Recursive Bayesian density learning

$$\mathcal{X}^n = x_{1:n} : p(\theta|\mathcal{X}^n) = \frac{p(x_n|\theta)p(\theta|\mathcal{X}^{n-1})}{\int p(x_n|\theta)p(\theta|\mathcal{X}^{n-1})d\theta}$$

Frequentist vs Bayesian

Bayes: priors, distributions, needs efficient integration, adds regularization term.

Frequentist: no priors, point estimate, requires only differentiation methods.

MLE are consistent, equivariant, asymptotically normal, asymptotically efficient (no efficient for finite samples).

Data Types

$$\text{monadic: } X: O \rightarrow \mathbb{R}^d \text{ dyadic: } X: O_1 \times O_2 \rightarrow \mathbb{R}^d.$$

$$\text{pairwise: } X: O_1 \times O_1 \rightarrow \mathbb{R}^d \text{ polyadic}$$

$$\text{data: } X: O_1 \times O_2 \times O_3 \rightarrow \mathbb{R}^d \text{ nominal =}$$

$$\text{qualitative (sweet, sour ...), ordinal =}$$

$$\text{absolute order, quantitative = numbers}$$

Regression

$$\text{Model of data: } \mathbf{Y} = \mathbf{X}\beta^* + \varepsilon$$

$$\mathbf{X} \in \mathbb{R}^{(d+1) \times n} \quad \beta \in \mathbb{R}^{d+1} \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbb{I}\sigma^2)$$

$$\mathbf{Y}|\mathbf{X}, \beta, \sigma^2 \sim \mathcal{N}(\mathbf{Y}; \mathbf{X}^T \beta, \mathbb{I}_{(d+1)} \sigma^2)$$

MLE: Ordinary Least Squares

OLSE is unbiased, orthogonal projection with lowest variance. differentiate wrt β .

$$\mathcal{L} = \text{RSS}(\beta) = \sum_{i=1}^n (y_i - x_i^T \beta)^2 = (\mathbf{y} - \mathbf{X}\beta)^2$$

$$\text{Estimator: } \hat{\beta}^{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$\text{Prediction: } \hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

MAP: Ridge Regression (L^2 penalty)

$$\text{Penalize energy in } \beta. \text{ Prior: } \beta \sim \mathcal{N}(\mathbf{0}, \lambda^{-1} \mathbb{I})$$

$$\text{Loss: } \mathcal{L} = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \beta^T \beta$$

$$\text{Estimator: } \hat{\beta}^{\text{ridge}} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbb{I})^{-1} \mathbf{X}^T \mathbf{y}$$

GDM: MAP: Lasso (L^1 penalty)

Penalize full β . Lasso has no closed form.

$$\beta \sim \text{Lapl}(0, \lambda^{-1}) = \frac{\lambda}{2} \exp(-\lambda|\beta|)$$

$$\mathcal{L} = \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \sum_{j=1}^d |\beta_j|$$

$$= (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) + \lambda \|\beta\|_1$$

$$\text{Bayesian view: } Y|(\mathbf{X}, \beta) \sim \mathcal{N}(x^T \beta, \sigma^2 \mathbb{I})$$

d-Dim Bayesian Linear Regression

$$\text{Prior: } \beta \sim \mathcal{N}(\mu_0, \Lambda^{-1})$$

$$\text{Likelihood: } Y|\beta, \mathbf{X}, \sigma \sim \mathcal{N}(X\beta, \sigma_n^2 \mathbb{I})$$

$$\text{Posterior: } \beta|\mathbf{X}, \mathbf{y} \sim \mathcal{N}(\mu, \Sigma)$$

$$\cdot \Sigma = (\sigma_n^{-2} \mathbf{X}^T \mathbf{X} + \Lambda)^{-1}$$

$$\cdot \mu = \Sigma(\Lambda \mu_0 + \sigma_n^{-2} \mathbf{X}^T \mathbf{y})$$

Nonlinear Regression

Idea: Add feature space transformation, kernel to compute inner product. Suppose:

$$\beta \sim \mathcal{N}(\mathbf{0}, \Lambda^{-1}) \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma_n^2 \mathbb{I}_d)$$

$$\mathbf{Y} = \mathbf{X}\beta + \varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{X}\Lambda^{-1} \mathbf{X}^T + \sigma_n^2 \mathbb{I}_d)$$

Kernels

$$\text{Kernel: } k(x_i, x_j) = \phi(x_i) \Lambda^{-1} \phi(x_j)^T$$

Similarity based reasoning.

$$\text{Gram Matrix: } K = k(\mathbf{x}_i, \mathbf{x}_j), \quad 1 \leq i, j \leq n$$

$$\cdot k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x}) \cdot k(\mathbf{x}, \mathbf{x}') \text{ pos. semi-def.}$$

$$\text{If } k_1, k_2 \text{ kernels, } c \in \mathbb{R}_{>0}, \mathbf{A}^{psd}, p_{\text{pos-coeff}}:$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}'$$

$$= k_1(\mathbf{x}, \mathbf{x}') \cdot k_2(\mathbf{x}, \mathbf{x}') = c \cdot k_1(\mathbf{x}, \mathbf{x}')$$

$$= p(k_1(\mathbf{x}, \mathbf{x}')) = f(\mathbf{x}) k_1(\mathbf{x}, \mathbf{x}') f(\mathbf{x}')$$

$$k(\mathbf{x}, \mathbf{x}') = \phi(x)^T \phi(x') = (1 + \mathbf{x}^T \mathbf{x}')^m$$

$$= \tanh(\alpha \mathbf{x}^T \mathbf{x}' + c)$$

$$= \sigma^2 \exp\left(-\frac{2 \sin^{-1} \pi \|\mathbf{x} - \mathbf{x}'\|_2^2}{l^2}\right)$$

$$= \exp(-\|\mathbf{x} - \mathbf{x}'\|_1 l^{-1})$$

$$= \exp(-\|\mathbf{x} - \mathbf{x}'\|_2^2 (2l^2)^{-1})$$

$$\text{RBF: } \phi_j(x) = \exp\left(-\frac{\|x\|_2^2}{2}\right) \prod_{i=0}^d x^{j_i} (j_i!)^{-\frac{1}{2}}$$

↑ Lengthscale, smoother fcts.

Gaussian Process Regression

Applying a kernel, we get:

$$\mathbf{Y} = \Phi \beta + \varepsilon \sim \mathcal{N}(\mathbf{0}, \Phi \Lambda^{-1} \Phi^T + \sigma_n^2 \mathbb{I}_d) =$$

$$\mathcal{N}\left(\begin{bmatrix} \mathbf{y} \\ \mathbf{y}_* \end{bmatrix} \middle| \mathbf{0}, \begin{bmatrix} \mathbf{K} + \sigma^2 \mathbb{I} & \mathbf{k} \\ \mathbf{k}^T & k(x_*, x_*) + \sigma^2 \end{bmatrix}\right)$$

Gaussian Process Prediction

$$\text{Given } \mathcal{GP}(\mu, K),$$

$$p(y_*|x_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(\tilde{\mu}, \tilde{\sigma}^2),$$

$$\cdot \tilde{\mu} = \mu(x_*) + \mathbf{k}^T (\mathbf{K} + \sigma_n^2 \mathbb{I})^{-1} (\mathbf{y} - \mu(\mathbf{X})),$$

$$\cdot \tilde{\sigma}^2 = k(x_*, x_*) - \mathbf{k}^T (\mathbf{K} + \sigma^2 \mathbb{I})^{-1} \mathbf{k}$$

$$\cdot \mathbf{k} = k(x_*, \mathbf{X}) \quad \mathbf{K}_{ij} = k(x_i, x_j)$$

$$\cdot \tilde{\sigma}_{ij}^2 = k(x_i, x_j) - \mathbf{k}_i^T (\mathbf{K} + \sigma^2 \mathbb{I})^{-1} \mathbf{k}_j$$

Causality

KCV-Cross Validation

Partition data Z into K equally sized, disjoint subsets:

$$\mathcal{Z} = \mathcal{Z}_1 \cup \mathcal{Z}_2 \cup \dots \cup \mathcal{Z}_K, \mathcal{Z}_\mu \cap \mathcal{Z}_\nu = \emptyset$$

$|\mathcal{Z}_k| \approx n \frac{K-1}{K}$ #of training samples. Learn:

$$\hat{f}^{-\nu}(x) = \arg \min_{f \in \mathcal{F}} \frac{\sum_{i \in \mathcal{Z}_\nu} \mathcal{L}(y_i, f(x_i))}{|\mathcal{Z} - \mathcal{Z}_\nu|}$$

$$\hat{R}^{CV}(A) = \frac{1}{n} \sum_{i \leq n} \mathcal{L}(y_i, \hat{f}^{-\kappa(i)}(x_i))$$

Underfits because smaller dataset.

Leave-one-out: $K = n$ (unbiased but var can be large from correlated datasets)

Bootstrapping

Bootstrap samples: $\mathcal{Z}^* = \{\mathcal{Z}_1^*, \dots, \mathcal{Z}_B^*\}$, of same size as original, drawn with replacement. The chance of a sample to have appeared in the bootstrap is:

$1 - (1 - \frac{1}{n})^n \xrightarrow{n \rightarrow \infty} 1 - \frac{1}{e} \approx 0.632$. So if we compute the ERM on \mathcal{Z} we could get 63% accuracy by memorization. Over-confident (shows too small bias)!

Leave-one-out/out of bucket error: compensates by computing the ERM where no memorization was for specific sample. E.g., for classification, like cross-validation:

$$\hat{R}(\mathcal{A}) = \frac{1}{B} \sum_{b=1}^B \sum_{z_i \in \mathcal{Z}^{*b}} \frac{\mathbb{I}_{c(z_i) \neq y_i}}{B - |\mathcal{Z}^{*b}|} \hat{R}_{0.632} = 0.368 \hat{R}(\mathcal{A}) + 0.632 \hat{R}_{bs}$$

$$\text{Wald Test: } W = \frac{\hat{\theta} - \theta_0}{\text{s.e.}(\hat{\theta})}$$

Bayesian Neural Networks (BNN)

NN: no uncertainty quantification, overconfident, adversarial examples, poor generalization for domain shifts.

BNN: Using $p(w)$ and $p(D|w)$, approx. poster. by variational infer. (min rev KL).

$$\sigma \leftarrow \sigma - \alpha_t \left(\varepsilon^T \frac{\partial}{\partial w} F(w, \theta) + \frac{\partial}{\partial \sigma} F(w, \theta) \right)$$

Information-based Transductive Lear.

ITL selects x_n that maximizes mutual information of $y_x = f_x + \varepsilon_x$ about f :

$$x_n = \arg \max_{x \in \mathcal{X}} I(f_A; y_x | D_{n-1})$$

If $f \sim \text{GP}(\mu, k)$, then:

$$I(f_A; y_x | D_{n-1}) = \frac{1}{2} \log \left(\frac{\text{Var}[y_x | D_{n-1}]}{\text{Var}[y_x | f_A, D_{n-1}]} \right)$$

Safe Bayesian Optimization

$$x_n = \arg \max_{x \in \mathcal{S}_n = \{x | u_n^g(x) \geq 0\}} u_n^f(x)$$

Batch Active Learning | ProbCover

$$G = (X, E), \quad E = \{(x, x') \mid \|x - x'\| \leq \delta\}$$

$$L \leftarrow \emptyset \quad \forall i = 1, 2, \dots, b \quad \hat{x} \leftarrow$$

$$\arg \max_{x \in X} |\{x' \mid (x, x') \in E, x' \in X\}|$$

$$L \leftarrow L \cup \{\hat{x}\} \mid E \leftarrow E \setminus (\{\hat{x}\} \times (B_\delta(\hat{x}) \cap X))$$

Classification

Definitions and Lemmas from Chapter 7

Definition 11 (Constrained Optimization Problem) A constrained

optimization problem is of the form:

$$\min_{w \in \mathbb{R}^d} f(w) \quad \text{s.t.} \quad g_i(w) = 0, i \leq m, \quad h_j(w) \leq 0, j \leq n$$

A feasible solution satisfies all constraints, and an optimal solution minimizes $f(w)$ over all feasible solutions.

Definition 12 (Convex Optimization Problem) A constrained optimization problem is convex if $f, g_1, \dots, g_m, h_1, \dots, h_n$ are convex and the feasible region is convex.

Definition 13 (Lagrangian) The Lagrangian of a constrained optimization problem is:

$$L(\lambda, \alpha, w) = f(w) + \sum_{i \leq m} \lambda_i g_i(w) + \sum_{j \leq n} \alpha_j h_j(w)$$

where λ and α are the Lagrange multipliers. **Lemma 14 (Necessary Conditions for Optimal Solutions)** Any optimal solution w^* must satisfy:

$$\nabla_w L(\lambda, \alpha, w) = 0, \quad g_i(w) = 0, \quad h_j(w) \leq 0$$

for $i \leq m, j \leq n$. **Lemma 15 (Dual Function and Lower Bound)** Given the Lagrangian $L(\lambda, \alpha, w)$, the dual function is:

$$\theta(\lambda, \alpha) = \inf_w L(\lambda, \alpha, w),$$

which satisfies:

$$\max_{\lambda, \alpha \geq 0} \inf_w L(\lambda, \alpha, w) \leq f(w^*).$$

This forms the basis of the dual problem. **Definition 17 (Strong Duality)** A convex optimization problem satisfies strong duality if:

$$\theta(\lambda^*, \alpha^*) = f(w^*),$$

where w^* solves the primal and (λ^*, α^*) solves the dual. **Definition 18 (Slater's Condition)** A convex optimization problem satisfies Slater's condition if there exists a strictly feasible point w_0 such that:

$$h_j(w_0) < 0, \quad \forall j \leq n.$$

This ensures strong

duality. **Lemma 19 (Slater's Condition and Strong Duality)** If a convex optimization problem satisfies Slater's condition, then strong duality holds.

Lemma 20 (Complementary Slackness) If strong duality holds, then for any optimal solution w^* :

$$f(w^*) = L(\lambda^*, \alpha^*, w^*),$$

$$\alpha_j h_j(w^*) = 0, \quad \forall j \leq n.$$

which enforces that inactive constraints have zero dual multipliers. **Bayes Optimal Classifier** Minimizes total risk for 0-1 Loss

$$\hat{c}(x) = \begin{cases} y & \mathbb{P}[y|x] > 1 - d, \exists y \\ \mathcal{D} & \mathbb{P}[y|x] < 1 - d, \forall y \end{cases}$$

Linear Classifier

$$g(x) = a^T \tilde{x} \quad a = (w_0, w)^T, \tilde{x} = (1, x)^T$$
$$a^T \tilde{x}_i > 0 \Rightarrow y_i = 1, a^T \tilde{x}_i < 0 \Rightarrow y_i = 2$$

Normalization: $\tilde{x}_i \rightarrow -\tilde{x}_i$ if $y_i = 2$

Find a : $a^T \tilde{x}_i > 0, \forall i$!

Perceptron Criterion

$$J_P(a) = \sum_{\tilde{x} \in \mathcal{X}^{\text{misl}}} (-a^T \tilde{x}), \quad \Rightarrow a_{k+1} = a_k + \eta_k \sum_{\tilde{x} \in \mathcal{X}^{\text{misl}}} \text{class} * \tilde{x}$$

Converges if data separable.

WINNOW Algorithm

Performs better when many dimensions are irrelevant. Search for 2 weight vectors a^+, a^- (for each class). If a point is misclassified: $a_i^+ \leftarrow \alpha^{-\tilde{x}_i} a_i^+, a_i^- \leftarrow \alpha^{-\tilde{x}_i} a_i^-$ (class 1 err.) $a_i^+ \leftarrow \alpha^{-\tilde{x}_i} a_i^+, a_i^- \leftarrow \alpha^{+\tilde{x}_i} a_i^-$ (class 2 err.) Exponential update.

Support Vector Machine (SVM)

Generalize Perceptron with margin and kernel. Find plane that max. margin m s.t.:

$$z_i g(\mathbf{y}) = z_i (\mathbf{w}^T \mathbf{y} + w_0) \geq m, \forall \mathbf{y}_i \in \mathcal{Y}$$

$$z_i \in \{-1, +1\} \quad \mathbf{y}_i = \phi(\mathbf{x}_i)$$

Vectors \mathbf{y}_i are the support vectors

Functional Margin Problem: minimizes $\|\mathbf{w}\|$ for $m=1$: $L(\mathbf{w}, w_0, \alpha) =$

$$= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^n \alpha_i [z_i (\mathbf{w}^T \mathbf{y}_i + w_0) - 1]$$

where α are Lagrange multipliers. $\frac{\partial L}{\partial \mathbf{w}} = 0$ and $\frac{\partial L}{\partial w_0} = 0$ give us constraints

$$\mathbf{w} = \sum_{i=1}^n \alpha_i z_i \mathbf{y}_i \quad 0 = \sum_{i=1}^n \alpha_i z_i$$

Replacing these in L we get (max α)

$$\tilde{L}(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j z_i z_j \mathbf{y}_i^T \mathbf{y}_j$$

with $\alpha_i \geq 0$ and $\sum_{i=1}^n \alpha_i z_i = 0$

This is the dual representation. The optimal hyperplane is given by

$$\mathbf{w} = \sum_{i=1}^n \alpha_i^* z_i \mathbf{y}_i$$

$$w_0^* = -\frac{1}{2} (\min_{z_i=1} \mathbf{w}^{*T} \mathbf{y}_i + \max_{z_i=-1} \mathbf{w}^{*T} \mathbf{y}_i)$$

where α maximize the dual problem.

Only Support Vectors ($\alpha_i \neq 0$) contribute to the evaluation.

$$\text{Optimal Margin: } \mathbf{w}^T \mathbf{w} = \sum_{i \in \text{SV}} \alpha_i^*$$

$$\text{Discrim.: } g^*(\mathbf{x}) = \sum_{i \in \text{SV}} z_i \alpha_i \mathbf{y}_i^T \mathbf{y} + w_0^*$$

$$\text{class} = \text{sign}(\mathbf{y}^T \mathbf{w}^* + w_0^*)$$

Kuhn-Tucker Conditions: only then strong duality holds: $\alpha_i^* \geq 0$

$$\alpha_i^* (z_i g^*(y_i) - 1) = 0, (z_i g^*(y_i) - 1) \geq 0$$

Soft Margin SVM

Introduce slack to relax constraints minimize $\frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i$ with respect

$$\text{to: } z_i (\mathbf{w}^T \mathbf{y}_i + w_0) \geq m(1 - \xi_i)$$

$$\text{Lagrangian: } L(\mathbf{w}, w_0, \xi, \alpha, \beta) = \frac{1}{2} \mathbf{w}^T \mathbf{w} +$$

$$C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i [z_i (\mathbf{w}^T \mathbf{y}_i + w_0) - 1 + \xi_i]$$

$$- \sum_{i=1}^n \beta_i \xi_i$$

C controls margin maximization vs. constraint violation

Dual Problem same as usual SVM but with supplementary constraint: $C \geq \alpha_i \geq 0$

Kuhn-Tucker Conditions: $\alpha_i^* (z_i (\mathbf{w}^T \mathbf{y}_i + w_0) - 1 + \xi) = 0, \xi_i (\alpha_i - C) = 0$

Non-Linear SVM

Use kernel in discriminant function:

$$g(\mathbf{x}) = \sum_{i,j=1}^n \alpha_i z_i K(\mathbf{x}_i, \mathbf{x})$$

E.g solve the XOR Problem with:

$$K(x, y) = (1 + x_1 y_1 + x_2 y_2)^2$$

Multiclass SVM

$\forall \text{class } z \in \{1, 2, \dots, M\}$ we introduce w_z and define our problem:

$$\min_w \frac{1}{2} \mathbf{w}^T \mathbf{w} = \min_{\{w_z\}_{z=1}^M} \sum_{z=1}^M \mathbf{w}_z^T \mathbf{w}_z$$

s.t. $(\mathbf{w}_z^T \mathbf{y}_i + w_{z,0}) -$

$$\max_{z \neq Z_i} (\mathbf{w}_z^T \mathbf{y}_i + w_{z,0}) \geq 1, \forall \mathbf{y}_i \in \mathcal{Y}$$

classification: $\hat{z} = \arg \max_z (\mathbf{w}_z^T \mathbf{y} + w_{z,0})$

Structured SVM

Each sample \mathbf{y} is assigned to a structured output label z

Output Space Representation:

joint feature map: $\psi(z, \mathbf{y})$

Scoring function: $f_w(z, \mathbf{y}) = \mathbf{w}^T \psi(z, \mathbf{y})$

Classify: $\hat{z} = h(\mathbf{y}) = \arg \max_{z \in \mathbb{K}} f_w(z, \mathbf{y})$

SVM objective:

$$\mathbf{w}^T \psi(z_i, \mathbf{y}_i) - \max_{z_i \neq z} \mathbf{w}^T \psi(z, \mathbf{y}_i) \geq m$$

with margin rescaling: $\min_{w, \xi \geq 0} \frac{1}{2} \mathbf{w}^T \mathbf{w} +$

$$C \sum_{i=1}^n \xi_i \quad \text{s.t.} \quad \mathbf{w}^T \psi(z_i, \mathbf{y}_i) - \Delta(z, z_i) -$$

$$\mathbf{w}^T \psi(z, \mathbf{y}_i) \geq -\xi_i, \forall z \neq z_i \quad \forall i$$

Lagrangian: let $\mathbb{K}_i = \mathbb{K} \setminus z_i$

$$\frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^n \xi_i -$$

$$\sum_{i=1}^n \sum_{z_j \in \mathbb{K}_i} \alpha_{i,j} (\mathbf{w}^T \psi(z_i, \mathbf{y}_i) -$$

$$\Delta(z_j, z_i) - \mathbf{w}^T \psi(z_j, \mathbf{y}_i) + \xi_i) -$$

$$\sum_{i=1}^n \beta_i \xi_i \quad \text{with } \alpha_{i,j} \geq 0, \beta_i \geq 0$$

Ensemble Methods

Combining Regressors

Set of estimators: $\hat{f}_1(x), \dots, \hat{f}_B(x)$

simple average: $\hat{f}(x) = \frac{1}{B} \sum_{i=1}^B \hat{f}_i(x)$

$$\text{Bias}[\hat{f}(x)] = \frac{1}{B} \sum_{i=1}^B \text{Bias}[\hat{f}_i(x)]$$

$$\mathbb{V}[\hat{f}(x)] \approx \frac{\sigma^2}{B} \quad \text{if the estimators are uncorrelated.}$$

Combining Classifiers

Input: classifiers $c_1(x), \dots, c_B(x)$

$$\text{Infer } \hat{c}_B(x) = \text{sgn}(\sum_{b=1}^B \alpha_b c_b(x))$$

with weights $\{\alpha_b\}_{b=1}^B$

Requires diversity of the classifiers.

Bagging

Train on bootstrapped subsets.

Sample: $\mathcal{Z} = \{(x_1, y_1), \dots, (x_n, y_n)\}$

\mathcal{Z}^* : chose i.i.d from \mathcal{Z} w. replacement.

Covariance small, variance similar, bias weakly affected.

Random Forest (Bagging strategy)

Collection of uncorr. decision trees.

Partition data space recursively. Grow the tree sufficiently deep to reduce bias. (each tree on other bagged set and with random collection of features available at every node)

Prediction with voting.

Boosting

Combine uncorr. weak learners in sequence. (Weak to avoid overfitting).

Coeff. of \hat{c}_{b+1} depend on \hat{c}_b 's results

AdaBoost (minimizes exp. loss)

Init: $\mathcal{X} = \{(x_1, y_1), \dots, (x_n, y_n)\}, w_i^{(1)} = \frac{1}{n}$

Fit $\hat{c}_b(x)$ to \mathcal{X} weighted by $w^{(b)}$

$$\varepsilon_b = \sum_{i=1}^n w_i^{(b)} \mathbb{I}_{\{\hat{c}_b(x_i) \neq y_i\}} / \sum_{i=1}^n w_i^{(b)}$$

$$\alpha_b = \log \frac{1 - \varepsilon_b}{\varepsilon_b} > 0$$

$$w_i^{(b+1)} = w_i^{(b)} \exp(\alpha_b \mathbb{I}_{\{\hat{c}_b(x_i) \neq y_i\}})$$

return $\hat{c}_B(x) = \text{sgn}(\sum_{b=1}^B \alpha_b \hat{c}_b(x))$

Best approx. at log-odds ratio.

Like stagewise-additive modeling.

Difference

(1) Boosting keeps identical training data,

bagging potentially varies the training data for each classifier. (2) Boosting weighs the prediction of each classifier according to its accuracy, bagging gives same importance to each.

Notes

AdaBoost gives large weight to samples that are hard to classify: those could be outliers. For bagging, there is a chance that imbalanced data-sets lead to bootstrap samples missing a class altogether. Fix by making the bootstrap size large enough s.t. at least one point is included.

$$\log \frac{P(y=1|x)}{P(y=-1|x)} = \sum_{b=1}^B c_b(x) =: F(x)$$

remove it and decrease K .