

1 Probability

Cond. Ind. $X \perp Y | Z \implies P(X, Y | Z) = P(X | Z)P(Y | Z)$

Cond. Ind. $X \perp Y | Z \implies P(X | Y, Z) = P(X | Z)$

$$\mathbb{E}[X] = \int_{\mathcal{X}} t \cdot f_X(t) dt =: \mu_X$$

$$\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \int_{\mathcal{X}} (t - \mathbb{E}[X])^2 f_X(t) dt = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

$$\text{Cov}(X, Y) = \int_{\mathcal{X}} \int_{\mathcal{Y}} f(x, y)(x - \mu_x)(y - \mu_y) dx dy$$

“ $\mathbf{X}^2 = \mathbf{X}\mathbf{X}^T \geq 0$ ((symmetric) positive semidefinite)

$$\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

$$\text{Var}[\mathbf{A}\mathbf{X}] = \mathbf{A} \text{Var}[\mathbf{X}] \mathbf{A}^T \quad \text{Var}[aX + b] = a^2 \text{Var}[X]$$

$$\text{Var}[\sum_{i=1}^n a_i X_i] = \sum_{i=1}^n a_i^2 \text{Var}[X_i] + 2 \sum_{i,j,i < j} a_i a_j \text{Cov}(X_i, X_j)$$

$$\text{Var}[\sum_{i=1}^n a_i X_i] = \sum_{i=1}^n a_i^2 \text{Var}[X_i] + \sum_{i,j,i \neq j} a_i a_j \text{Cov}(X_i, X_j)$$

$$\frac{\partial}{\partial t} P(X \leq t) = \frac{\partial}{\partial t} F_X(t) = f_X(t) \text{ (derivative of c.d.f. is p.d.f)}$$

$$f_{\alpha Y}(z) = \frac{1}{\alpha} f_Y(\frac{z}{\alpha})$$

$$\text{Empirical CDF: } \hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\{X_i \leq t\}}$$

$$\text{Empirical PDF: } \hat{f}_n(t) = \frac{1}{n} \sum_{i=1}^n \delta(t - X_i) \text{ (continuous)}$$

$$\text{Empirical PDF: } \hat{p}_n(t) = \frac{1}{n} |x = t| x \in D \text{ (discrete)}$$

T. The MGF $\psi_X(t) = \mathbb{E}[e^{tX}]$ characterizes the distr. of a rv

$$Be(p): \quad pe^t + (1-p) \quad \left| \quad \mathcal{N}(\mu, \sigma): \quad \exp\left(\mu t + \frac{1}{2}\sigma^2 t^2\right)$$

$$Bin(n, p): (pe^t + (1-p))^n \quad \left| \quad Gam(\alpha, \beta): \left(\frac{1}{a-\beta t}\right)^\alpha \text{ for } t < 1/\beta$$

$$Pois(\lambda): e^{\lambda(e^t-1)}$$

T. If X_1, \dots, X_n are ind. rvs with MGFs $M_{X_i}(t) = \mathbb{E}[e^{tX_i}]$, then the MGF of $Y = \sum_{i=1}^n a_i X_i$ is $M_Y(t) = \prod_{i=1}^n M_{X_i}(a_i t)$.

T. Let X, Y be ind., then the p.d.f. of $Z = X + Y$ is the conv. of the p.d.f. of X and Y :

$$f_Z(z) = \int_{\mathbb{R}} f_X(t) f_Y(z-t) dt = \int_{\mathbb{R}} f_X(z-t) f_Y(t) dt$$

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

$$\mathbf{T}. P\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix}\right) = \mathcal{N}\left(\begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \end{bmatrix} \middle| \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{u}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}\right)$$

$$\mathbf{a}_1, \mathbf{u}_1 \in \mathbb{R}^e, \boldsymbol{\Sigma}_{11} \in \mathbb{R}^{e \times e} \text{ p.s.d.}, \boldsymbol{\Sigma}_{12} \in \mathbb{R}^{e \times f} \text{ p.s.d.}$$

$$\mathbf{a}_2, \mathbf{u}_2 \in \mathbb{R}^f, \boldsymbol{\Sigma}_{22} \in \mathbb{R}^{f \times f} \text{ p.s.d.}, \boldsymbol{\Sigma}_{21} \in \mathbb{R}^{f \times e} \text{ p.s.d.}$$

$$P(\mathbf{a}_2 | \mathbf{a}_1 = \mathbf{z}) = \mathcal{N}(\mathbf{a}_2 | \mathbf{u}_2 + \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1}(\mathbf{z} - \mathbf{u}_1), \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12})$$

T. (Chebyshev) Let X be a rv with $\mathbb{E}_X [=] \mu$ and variance $\text{Var}[X] = \sigma^2 < \infty$. Then for any $\epsilon > 0$, we have $P(|X - \mu| \geq \epsilon) \leq \frac{\sigma^2}{\epsilon^2}$.

T. (Cramer Rao Bound) If $\hat{\theta}$ is unb., cons. est., then

$$MSE(\hat{\theta}) = \mathbb{E}\left[\left(\hat{\theta} - \theta\right)^2\right] \geq \frac{1}{\mathcal{I}_n(\theta)} > 0, \quad \text{where}$$

$$\mathcal{I}_n(\theta) = -\sum_{i=1}^n \mathbb{E}\left[\frac{\partial^2 \log(P(x_i | \boldsymbol{\theta}))}{\partial \theta^2}\right]. \quad \text{(Fisher Information)}$$

D. (Conditional Expected Risk) Given rv X

$$R(f, X) = \int_{\mathcal{Y}} L(Y, f(X)) P(Y | X) dY$$

D. (Total Expected Risk) for rvs X, Y

$$R(f) = \mathbb{E}_X [R(f, X)] = \int_{\mathcal{X}} R(f, X) P(X) dX = \mathbb{E}_{X,Y} [L(Y, f(X))].$$

2 Matrix Calculus

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \nabla_{\mathbf{x}} f(\mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) \quad \text{(2nd order Taylor Expan.)}$$
$$+ \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^T \text{Hess}(f; \mathbf{x}_0) (\mathbf{x} - \mathbf{x}_0) + \mathcal{O}(\|\mathbf{x} - \mathbf{x}_0\|^3)$$

$$\text{Hess}(f; \mathbf{x}_0) = \frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{x}^T}(\mathbf{x}_0)$$

$$\frac{\partial}{\partial \mathbf{x}} [\mathbf{A}\mathbf{x}] = \mathbf{A}^T \quad \frac{\partial}{\partial \mathbf{x}} [\mathbf{A}f(\mathbf{x})] = \mathbf{A}^T \frac{\partial}{\partial \mathbf{x}} [f(\mathbf{x})]$$

$$\frac{\partial}{\partial \mathbf{x}} [\mathbf{x}^T \mathbf{A}\mathbf{x}] = 2\mathbf{A}^T \mathbf{x} \quad \frac{\partial}{\partial \mathbf{x}} \left[\|\mathbf{f}(\mathbf{x})\|_2^2\right] = 2 \frac{\partial}{\partial \mathbf{x}} [\mathbf{f}(\mathbf{x})] \mathbf{f}(\mathbf{x})$$

T. (Sylvester Criterion) A $d \times d$ matrix is positive semi-definite if and only if all the upper left $k \times k$ for $k = 1, \dots, d$ have a positive determinant.

3 Misc

T. (Jensen) f convex/concave, $\forall i: \lambda_i \geq 0, \sum_{i=1}^n \lambda_i = 1$

$$f\left(\sum_{i=1}^n \lambda_i \mathbf{x}_i\right) \leq / \geq \sum_{i=1}^n \lambda_i f(\mathbf{x}_i)$$

Special case: $f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$.

T. (CSU) $|\langle \mathbf{x}, \mathbf{y} \rangle|^2 \leq \langle \mathbf{x}, \mathbf{x} \rangle \cdot \langle \mathbf{y}, \mathbf{y} \rangle$

$$\text{Special case: } \left(\sum x_i y_i\right)^2 \leq \left(\sum x_i^2\right) \left(\sum y_i^2\right)$$

$$\text{Special case: } \mathbb{E}[XY]^2 \leq \mathbb{E}[X^2] \mathbb{E}[Y^2]$$

Lag.: $f(x, y)$ s.t. $g(x, y) = c \quad \mathcal{L}(x, y, \gamma) = f(x, y) - \gamma(g(x, y) - c)$

4 Kernels

D. (Kernel) Kernel functions $k(x, x')$ must satisfy (cf. properties of covariance matrices)

1. Symmetry: $k(x, x') = k(x', x)$
2. Positive semi-definiteness (continuous case):
 $\forall f \in L_2 \forall \Omega \subset \mathbb{R}: \int_{\Omega} k(x, x') f(x) f(x') dx dx' \geq 0$
3. Positive semi-definiteness (discrete case):

For any $n \in \mathbb{N}$, and any set $S = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, the kernel (Gram) matrix $\mathbf{K} = (k(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^n$ must be positive semidefinite ($\forall \mathbf{x}: \mathbf{x}^T \mathbf{K} \mathbf{x} \geq 0$).

T. (Comp. Rules) $k_1 + k_2; k_1 \cdot k_2; c \cdot k_1 (c > 0); f(k_1)$ where f is a poly. w. pos. coeff., or the exp. func..

Kernel $k(\mathbf{x}, \mathbf{y}) =$	Feature Map $\phi(\mathbf{x}) =$
$k_1(\mathbf{x}, \mathbf{y}) + k_2(\mathbf{x}, \mathbf{y})$	$(\phi_1(\mathbf{x}), \phi_2(\mathbf{x}))^T$
$c \cdot k_1(\mathbf{x}, \mathbf{y}) (c > 0)$	$\sqrt{c} \cdot \phi_1(\mathbf{x})$
$k_1(\mathbf{x}, \mathbf{y}) \cdot k_2(\mathbf{x}, \mathbf{y})$	$(\phi_1(\mathbf{x})_i \cdot \phi_2(\mathbf{y})_j)^T$ for $1 \leq i \leq d_1, 1 \leq j \leq d_2$
$\mathbf{x}^T \mathbf{A} \mathbf{y} \quad \mathbf{A}$ p.s.d.	$\mathbf{L}^T \mathbf{x}, \mathbf{A} = \mathbf{L} \mathbf{L}^T$
$\mathbf{x}^T \mathbf{M}^T \mathbf{M} \mathbf{x}, \mathbf{M}$ arbitrary	$\mathbf{M} \mathbf{x}$

$$\text{Linear } k(x, y) = x^T y, \text{ Polynomial } k(x, y) = (x^T y + 1)^d$$

$$\text{RBF } k(x, y) = \exp\left(-\frac{\|x-y\|_2^2}{h^2}\right), \text{ Sigmoid } k(x, y) = \tanh(kx^T y - b)$$

Proof Tricks: Counterexample, Prove by “for arbitrary $n, x_1, \dots, x_n, c_1, \dots, c_n \sum_i \sum_j c_i c_j k(x_i, x_j) \geq 0$ ”, constant kernel

$k(x, y) = c$ is a kernel, find feature map and inner product.

5 Regression

Problem What is the optimal estimate of a function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ based on noisy data $y_i = f(x_i) + \epsilon_i$

Solution the regression function

$$f^*(x) = \mathbb{E}[Y | X = x] = \int_{\mathcal{Y}} y \cdot P(y | X = x) dy$$

5.1 — Ridge Regression

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \quad (\lambda > 0, \text{ chosen via CV})$$

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad \text{(always has a solution)}$$

$$\mathbf{g}_t = -2\mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w}_t) + 2\lambda \mathbf{w}_t$$

Note: Scale of the data matters for $\lambda!$ (\rightarrow normalize data)

$$Y \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}, \sigma^2), \quad y_i = \mathbf{w}^T \mathbf{x}_i + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$P(Y = y | \mathbf{x} = \mathbf{x}, \boldsymbol{\theta} = (\mathbf{w}, \sigma^2)) = \mathcal{N}(y; h(\mathbf{x}), \sigma^2), \quad h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

$$\text{Weights prior: } \mathbf{w} \sim \mathcal{N}(0, \beta^2 \mathbf{I}), \quad w_i \sim \mathcal{N}(0, \beta^2)$$

Maximizing $P(\mathbf{w} | D)$ then leads to the connection $\lambda = \frac{\sigma^2}{\beta^2}$.

5.2 — Bayesian Linear Regression

$$Y = \mathbf{X}^T \boldsymbol{\beta} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$P(Y | \mathbf{X}, \boldsymbol{\beta}, \sigma) = \mathcal{N}(\mathbf{X}^T \boldsymbol{\beta}, \sigma^2) \propto \exp\left(-\frac{1}{2\sigma^2} (Y - \mathbf{X}^T \boldsymbol{\beta})^2\right)$$

$$P(\boldsymbol{\beta} | \boldsymbol{\Lambda}) = \mathcal{N}(\mathbf{o}, \boldsymbol{\Lambda}^{-1}) \propto \exp\left(-\frac{1}{2} \boldsymbol{\beta}^T \boldsymbol{\Lambda} \boldsymbol{\beta}\right)$$

$$P(\boldsymbol{\beta} | \mathbf{X}, \mathbf{y}, \boldsymbol{\Lambda}, \sigma) = \frac{P(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\Lambda}, \sigma) P(\boldsymbol{\beta} | \mathbf{X}, \boldsymbol{\Lambda}, \sigma)}{P(\mathbf{y} | \mathbf{X}, \boldsymbol{\Lambda}, \sigma)} \propto P(\mathbf{y} | \mathbf{X}, \boldsymbol{\beta}, \sigma) P(\boldsymbol{\beta} | \boldsymbol{\Lambda})$$

$$= \prod_{i=1}^n P(y_i | \mathbf{x}_i, \boldsymbol{\beta}, \sigma) P(\boldsymbol{\beta} | \boldsymbol{\Lambda}) = \mathcal{N}(\boldsymbol{\beta}; \boldsymbol{\mu}_{\boldsymbol{\beta}}, \boldsymbol{\Sigma}_{\boldsymbol{\beta}})$$

$$\boldsymbol{\mu}_{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X} + \sigma^2 \boldsymbol{\Lambda})^{-1} \mathbf{X}^T \mathbf{y}, \quad \boldsymbol{\Sigma}_{\boldsymbol{\beta}} = \sigma^2 (\mathbf{X}^T \mathbf{X} + \sigma^2 \boldsymbol{\Lambda})^{-1}$$

Special Case: Ridge Regression: $\boldsymbol{\Lambda} = \lambda \mathbf{I}, \sigma = 1$

5.3 — Kernelized Ridge Regression

Insight optimal \mathbf{w}^* lies in the span of the data.

$$\mathbf{w}^* = \mathbf{X}_{\phi}^T \mathbf{z}^* \quad (\mathbf{K} = \mathbf{X}_{\phi} \mathbf{X}_{\phi}^T \in \mathbb{R}^{n \times n})$$

$$\mathbf{z}^* = \arg \min_{\mathbf{z}} \|\mathbf{K}\mathbf{z} - \mathbf{y}\|_2^2 + \lambda \mathbf{z}^T \mathbf{K} \mathbf{z}$$

$$\mathbf{1) Closed form } \mathbf{z}^* = (\mathbf{X}_{\phi} \mathbf{X}_{\phi}^T + \lambda \mathbf{I})^{-1} \mathbf{y} = (\mathbf{K} - \lambda \mathbf{I})^{-1} \mathbf{y}$$

$$\mathbf{2) Gradient descent } \mathbf{g}_t = 2\mathbf{K}^T (\mathbf{K}\mathbf{z} - \mathbf{y}) + 2\lambda \mathbf{K} \mathbf{z}$$

$$\text{Prediction } f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x}) = \dots = \sum_{i=1}^n z_i k(\mathbf{x}_i, \mathbf{x})$$

Bayesian Interpretation Same as ridge regression, except that the hypothesis class for \mathcal{H} for h (comp. of mean) may be different.

5.4 — Sparse Regression: LASSO

Requires $\|\mathbf{w}\|_1 \leq s$. Prior: $w_i \sim p(w_i; 0, b) = \frac{1}{2b} \exp\left(-\frac{|w_i - \mu|}{b}\right)$

where $\mu = 0$, (connection: $\lambda = \frac{2\sigma^2}{b}$).

5.5 — Ensemble of Regressions

T. The average of B unbiased estimators \hat{f}_i remains unbiased.

$$\mathbb{E}\left[\frac{1}{B} \sum_{i=1}^B \hat{f}_i(X)\right] - \theta = 0.$$

T. The variance of an average of B unbiased estimators, (*)

which have small covariances ≈ 0 , and similar variances $\approx \sigma^2$ my reduce the variance to σ^2/B .

$$\text{Var} \left[\frac{1}{B} \sum_{i=1}^B \hat{f}_i(X) \right] = \frac{1}{B^2} \sum_{i=1}^B \text{Var} [\hat{f}_i(X)] + \frac{1}{B^2} \sum_{\substack{i,j \\ i \neq j}}^B \text{Cov} (\hat{f}_i(X), \hat{f}_j(X)) \stackrel{(*)}{\approx} \frac{\sigma^2}{B}$$

5.6 – Bias Variance for Squared Loss

$$\begin{aligned} \mathbb{E}_D \left[\mathbb{E}_{\mathbf{X}, Y} \left[(Y - \hat{h}_D(\mathbf{X}))^2 \right] \right] & \quad (\text{Note: } h^*(X) = \mathbb{E}_Y [Y | X]) \\ = \mathbb{E}_{\mathbf{X}} \left[\underbrace{\left(\mathbb{E}_D [\hat{h}_D(\mathbf{X})] - h^*(\mathbf{X}) \right)^2}_{\text{Bias}^2} \right] & + \mathbb{E}_{\mathbf{X}} \left[\underbrace{\mathbb{E}_D \left[\left(\hat{h}_D(\mathbf{X}) - \mathbb{E}_D [\hat{h}_D(\mathbf{X})] \right)^2 \right]}_{\text{Variance}} \right] \\ + \mathbb{E}_{\mathbf{X}, Y} \left[(Y - h^*(\mathbf{X}))^2 \right] & \quad (\text{Noise}) \end{aligned}$$

Derivation: 1) $\pm \mathbb{E}_Y [Y | X]$ gives noise, 2) Prove vanishing cross-term, 3) $\pm \mathbb{E}_D [\hat{f}(X)]$ gives bias and variance

5.7 – Gaussian Processes

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n), \quad \boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Lambda}^{-1})$$

Since the outputs \mathbf{y} are a linear combination of normally distributed rvs $\boldsymbol{\beta}$, they are jointly Gaussian themselves.

$$\mathbb{E}_{\boldsymbol{\beta}, \boldsymbol{\epsilon}} [\mathbf{y}] = \mathbb{E}_{\boldsymbol{\beta}, \boldsymbol{\epsilon}} [\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}] = \mathbf{X}\mathbb{E}_{\boldsymbol{\beta}} [\boldsymbol{\beta}] + \mathbb{E}_{\boldsymbol{\epsilon}} [\boldsymbol{\epsilon}] = \mathbf{X}\mathbf{0} + \mathbf{0} = \mathbf{0}$$

$$\text{Cov}(\mathbf{y}) = \mathbb{E} [\mathbf{y}\mathbf{y}^\top] = \mathbb{E} [(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon})(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon})^\top] = \mathbf{X}\boldsymbol{\Lambda}^{-1}\mathbf{X}^\top + \sigma^2 \mathbf{I}$$

Special case: $\boldsymbol{\Lambda}^{-1} := \lambda \mathbf{I}$: $\text{Cov}(\mathbf{y}) = \lambda^{-1}(\mathbf{X}\mathbf{X}^\top + \lambda \sigma^2 \mathbf{I}_n)$.

Then we can rewrite the joint distr. as $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{K} + \sigma^2 \mathbf{I})$

where $k(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^\top \boldsymbol{\Lambda}^{-1} \mathbf{x}_j$ is a so-called kernel function (could be replaced with others).

$$\begin{aligned} \text{Joint D. } P((y_{n+1}^{\mathbf{y}} | x_{n+1}, \mathbf{X}, \sigma) &= \mathcal{N} \left((y_{n+1}^{\mathbf{y}} | \mathbf{0}, \begin{pmatrix} \mathbf{C}_p & \mathbf{k} \\ \mathbf{k}^\top & c \end{pmatrix} \right) \\ \mathbf{K} = k(\mathbf{X}, \mathbf{X}), \quad \mathbf{C}_n = \mathbf{K} + \sigma^2 \mathbf{I}, \quad c &= k(x_{n+1}, x_{n+1}) + \sigma^2, \\ \mathbf{k} = k(x_{n+1}, \mathbf{X}) \end{aligned}$$

$$\begin{aligned} \text{Pred D. } P(y_{n+1} | x_{n+1}, \mathbf{X}, \mathbf{y}, \sigma) &= \mathcal{N} \left(y_{n+1} | \boldsymbol{\mu}_{y_{n+1}}, \sigma_{y_{n+1}}^2 \right) \\ \boldsymbol{\mu}_{y_{n+1}} = \mathbf{k}^\top \mathbf{C}_n^{-1} \mathbf{y} = \mathbf{k}^\top (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, \quad \sigma_{y_{n+1}}^2 &= c - \mathbf{k}^\top \mathbf{C}_n^{-1} \mathbf{k} \end{aligned}$$

6 Classification

$$0/1 \text{ Loss } w^* = \arg \min_w \sum_{i=1}^n [y_i \neq \text{sign}(w^\top x_i)]$$

$$\text{Perceptron } w^* = \arg \min_w \sum_{i=1}^n [\max(0, y_i w^\top x_i)]$$

$$\text{Exponential Loss } L(y, z) = \exp(-(2y-1)(2z-1))$$

$$\text{Logistic Loss } L(y, z) = \ln(1 + \exp(-(2y-1)(2z-1)))$$

Dep. on appl. use $\mathbf{y}\mathbf{w}^\top \mathbf{x}$ or $-(2y-1)(2z-1)$ as error

6.1 – Generative VS Discriminative

$$P(Y | X) = \frac{P(X | Y) P(Y)}{\sum_y P(X | Y) P(Y)}$$

Generative: Est. both $P(Y)$ and $P(X | Y)$ then use Bayes.

Discr.: Est. $P(Y | X)$ directly, fitting discr. f. $g(Y, X)$.

6.2 – SVMs

6.2.1 – Primal, constrained

$$\min_w w^\top w + C \sum_{i=1}^n \xi_i, \quad \text{s.t. } y_i w^\top x_i \geq 1 - \xi_i, \xi_i \geq 0$$

6.2.2 – Primal, unconstrained

$$\min_w w^\top w + C \sum_{i=1}^n \max(0, 1 - y_i w^\top x_i) \quad (\text{hinge loss})$$

6.2.3 – Dual

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^\top x_j, \quad \text{s.t. } 0 \leq \alpha_i \leq C$$

6.2.4 – Dual to Primal

Dual to primal: $w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$, $\alpha_i > 0$: support vector.

6.3 – Kernelized SVMs

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j), \quad \text{s.t. } 0 \leq \alpha_i \leq C$$

Classify: $y = \text{sign}(\sum_{i=1}^n \alpha_i y_i k(x_i, x))$

– How to find a^T ?

$a = \{w_0, w\}$ used along $\tilde{x} = \{1, x\}$

Gradient Descent: $a(k+1) = a(k) - \eta(k) \nabla J(a(k))$

Newton: 2nd ord. Taylor, where $\eta_{opt} = H^{-1}$, $H = \frac{\partial^2 J}{\partial a_i \partial a_j}$

J is the cost mat., popular: Perceptron cost: $J_p(a) = \sum (-a^T \tilde{x})$

6.4 – Logistic Regression

$$P(Y = y | \mathbf{x}, \mathbf{w}) = \text{Ber}(y; \sigma(\mathbf{w}^\top \mathbf{x})) = \frac{1}{1 + e^{-\mathbf{y}\mathbf{w}^\top \mathbf{x}}} = p_y$$

$$P(Y = +1 | \mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x}) = p_+$$

Gr. st. w. Gau. Pr.: $\mathbf{w} \leftarrow \mathbf{w}(1 - 2\lambda \eta_t) + \eta_t \mathbf{y} \mathbf{x} \hat{P}(Y = -y | \mathbf{w}, \mathbf{x})$

6.4.1 – Multi-Class Logistic Regression

$$P(Y = i | \mathbf{x}, \mathbf{w}_1, \dots, \mathbf{w}_c) = \frac{\exp(\mathbf{w}_i^\top \mathbf{x})}{\sum_{j=1}^c \exp(\mathbf{w}_j^\top \mathbf{x})} = p_i$$

7 Multiclass Classification

One-VS-All $y = \arg \max_{i \in \{1, \dots, c\}} f_i(\mathbf{x})$

One-VS-One $\binom{c}{2}$ bin. clf. for each $(i, j) \in \{1, \dots, c\}^2$.

$$f_{(i,j)}: \mathcal{X} \rightarrow \{-1, +1\} \quad y = \arg \max_{i \in \{1, \dots, c\}} \sum_{j=1, j \neq i}^c \mathbb{1}_{\{f_{(i,j)}(\mathbf{x}) = +1\}}$$

8 Jackknife

Method to compensate for syst. est. errors (bias reduction).

Goal: Numerically estimate the bias of an estimator \hat{S}_n .

$$\tilde{S}_n = \frac{1}{n} \sum_{i=1}^n S_{n-1}^i \quad (\text{LOO-Estimator})$$

Then the est. for the bias of \hat{S}_n is: $\text{bias}^{JK} = (n-1)(\tilde{S}_n - \hat{S}_n)$

Then the unbiased estimate is $\hat{S}_n^{JK} = \hat{S}_n - \text{bias}^{JK}$.

9 Probabilistic Methods

$$P(\text{model } \theta | \text{data } D) = \underbrace{P(\text{data} | \text{model})}_{\text{Likelihood}} \times \underbrace{P(\text{model})}_{\text{Prior}} = \underbrace{P(\text{data})}_{\text{Evidence}}$$

9.0.1 – Maximum (Cond.) LH Est., (MLE)

$$\begin{aligned} \boldsymbol{\theta}^* &= \arg \max_{\boldsymbol{\theta}} \hat{P}(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n, \boldsymbol{\theta}) \\ &= \arg \min_{\boldsymbol{\theta}} - \sum_{i=1}^n \log \hat{P}(y_i | \mathbf{x}_i, \boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \dots \text{insert, derivate.} \end{aligned}$$

9.0.2 – Maximum a Posteriori Estimate, (MAP)

$$\begin{aligned} \boldsymbol{\theta}^* &= \arg \max_{\boldsymbol{\theta}} P(\boldsymbol{\theta} | D) = \arg \max_{\boldsymbol{\theta}} P(\boldsymbol{\theta} | \mathbf{x}_1, \dots, \mathbf{x}_n, y_1, \dots, y_n) \\ &\stackrel{\text{i.i.d.}}{=} \arg \min_{\boldsymbol{\theta}} - \log P(\boldsymbol{\theta}) - \sum_{i=1}^n \log P(y_i | \mathbf{x}_i, \boldsymbol{\theta}) = \dots \text{insert, derivate} \end{aligned}$$

10 Ensemble Methods

Use combination of simple hyp. (weak learners) that are sufficiently diverse to prod. a valid sol. with low bias and var.

Bagging: train weak l. on bootstr. sets with equal weights.

Boosting: train on all data, but reweigh misclassified samples higher and use error-sensitive reweighting of classifiers.

Left out: Decision trees (Stump: $h(x) = \text{sign}(ax_i - t)$), Random forest: bagging of trees.

10.0.1 – Ada Boost

$$f^*(x) = \arg \min_{f \in F} \sum_{i=1}^n \exp(-y_i f(x_i)) \quad (\text{eff. mins. exp. l.})$$

Algorithm 1: AdaBoost Algorithm

Initialize the observation weights: $\forall i: w_i^{(1)} \leftarrow \frac{1}{n}$

for $b \leftarrow 1$ **to** B **do**

Fit a clf. $c_b(x)$ to the weighed training data (acc. to w_i)

$$\epsilon_b \leftarrow \frac{\sum_{i=1}^n w_i^{(b)} \mathbb{1}_{\{c_b(x_i) \neq y_i\}}}{\sum_{i=1}^n w_i^{(b)}} \quad (\text{error of the weak learner})$$

$$\alpha_b \leftarrow \log \left(\frac{1 - \epsilon_b}{\epsilon_b} \right) \quad (\text{weigh classifier acc. to accuracy})$$

Update the datapoint weights for the next step:

For all $i \in \{1, \dots, n\}$ do

$$w_i^{(b+1)} \leftarrow w_i^{(b)} \exp(\alpha_b \mathbb{1}_{\{y_i \neq c_b(x_i)\}})$$

return $\hat{c}_B(x) = \text{sign} \left(\sum_{b=1}^B \alpha_b c_b(x) \right)$

Additive log. reg., Bayesian approach (assumes poster.), Newtonlike updates (GD), if prev. clf. bad, next has high weight.

11 Generative Methods

11.1 – Naive Bayes

Features indep. $P(y|x) = \frac{1}{Z} P(y) P(x|y)$, $Z = \sum_y P(y) P(x|y)$

$$y = \arg \max_{y'} P(y' | x) = \arg \max_{y'} \hat{P}(y') \prod_{i=1}^d \hat{P}(x_i | y')$$

Discr. Func.: $f(x) = \log \left(\frac{P(y=1|x)}{P(y=-1|x)} \right)$, $y = \text{sign}(f(x))$

11.2 – Fischer's LDA

$$J(\mathbf{w}) = \frac{\|\mathbf{w}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)\|}{\mathbf{w}^\top (\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2) \mathbf{w}} \quad \hat{\mathbf{w}} = (\boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2)^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \quad (\text{unscaled})$$

However, all samples influence boundary (better: points at border, SVM)

12 Unsupervised Learning

12.1 – Gaussian Mixture Modeling

$$(\mu^*, \Sigma^*, \pi^*) = \arg \min - \sum_i \log \sum_{j=1}^k \pi_j \mathcal{N}(x_i | \mu_j, \Sigma_j)$$

12.2 – EM Algorithm

Problem: sum within log-term of likelihood.

E-step: expectation: pick clusters for points.

$$\text{Calculate } \gamma_j^{(t)}(x_i) = \frac{P(c|\theta^j)(x_i|c, \theta^j)}{\sum P(x_i|\theta)}$$

M-Step: maximum LH: adjust clusters to best fit points.

$$\text{prior}_j^{(t)} = \pi_j^{(t)} \leftarrow \frac{1}{n} \sum_{i=1}^n \gamma_j^{(t)}(x_i)$$

$$\mu_j^{(t)} \leftarrow \frac{\sum_{i=1}^n \gamma_j^{(t)}(x_i)(x_i)}{\sum_{i=1}^n \gamma_j^{(t)}(x_i)}, \quad \Sigma_j^{(t)} \leftarrow \frac{\sum_{i=1}^n \gamma_j^{(t)}(x_i)(x_i - \mu_j^{(t)})(x_i - \mu_j^{(t)})^\top}{\sum_{i=1}^n \gamma_j^{(t)}(x_i)}$$