

## Representations.

Conditional expected risk:

$$R(f, X) = \int y Q(X, f(x)) P(X, Y) dX dY.$$

Total expected risk:

$$R(\hat{f}) = \int_X \int_Y Q(Y, f(x)) P(X, Y) dX dY.$$

Empirical risk:  $\hat{R}(\hat{f}) = \frac{1}{n} \sum_{i=1}^n Q(Y_i, \hat{f}(X_i))$

## Density Estimation in Regression

Bayes Rule:  $P(\theta|x) = \frac{p(x|\theta)p(\theta)}{p(x)}$

Cramer-Rao bound:  $\mathbb{E}[(\hat{\theta} - \theta_0)^2] = \frac{1}{I(\theta_0)}$

Fisher info:  $I(\theta_0) = n \mathbb{V}_{\theta}[\log \log p(x|\theta)]$

Efficiency:  $\mathbb{E}[(\hat{\theta} - \theta_0)^2] = \frac{1}{I(\theta_0)}$

Bias - Var - Tradeoff

Bias - Var - Tradeoff

Bias ↑ Var ↓  $\rightarrow$  overfit  $\rightarrow$  Var ↑ Bias ↓

Var ↑  $\rightarrow$  underfit  $\rightarrow$  Var ↓ Bias ↑

Tradeoff:  $(\mathbb{E}[Y|X=x] - \hat{Y})^2 + (\hat{Y} - \mathbb{E}[Y|X=x])^2$

+ ( $\mathbb{E}[\hat{Y}(x) - \mathbb{E}[Y|X=x]]$ )<sup>2</sup>

Conditional Gaussian dist.

Density Estimation in Regression

Summary on MLE:

① Consistency:  $\hat{\theta}_n \xrightarrow{P} \theta_0$

② Equivariance:  $g(\hat{\theta}_M)$  is MLE of  $g(\theta)$  (if  $g$  is invertible)

③ Asymp. normality:  $\sqrt{n}(\hat{\theta} - \theta_0) \rightsquigarrow N(0, J\hat{J}^{-1})$

$\sqrt{n}(\theta) = -\mathbb{E}\left[\frac{\partial^2 \log P(\theta)}{\partial \theta \partial \theta^T}\right]$  Inv. =  $\nabla \left[ \frac{\partial \log P(\theta)}{\partial \theta} \right]$  Bayes linear regression

④ Asymp. efficiency:  $\frac{n}{n-2\alpha} (V(\hat{\theta}_M|X_1, \dots, X_n))^2 \rightsquigarrow 1$

Assume.  $p(\mathbf{x}|Y) \sim N(\mu, \sigma^2)$ ,  $p(\mu) \sim N(\mu_0, \sigma_0^2)$

$p(\mu|Y) \sim N(\mu_n, \sigma_n^2)$ ,  $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n Y_i$

$\sigma_n^2 = \frac{n\sigma^2}{n\sigma^2 + \sigma^2} \hat{\sigma}_n^2 + \frac{\sigma^2}{n\sigma^2 + \sigma^2} \sigma_0^2$

$p(\mathbf{x}|Y) \sim N(\mu_n, \sigma_n^2)$

Multivariate case:

$\Sigma^{-1} = n\Sigma^{-1} + Z^{-1}$ ,  $\Sigma^{-1} \mu_0 = n\Sigma^{-1} \mu_0 + Z^{-1} \mu_0$

Recursive.  $p(\theta|X^n) = \frac{p(\mathbf{x}|\theta)p(\theta)}{\int p(\mathbf{x}|\theta)p(\theta) d\theta}$

## Regression.

$$\hat{\beta} \text{ RSS} = (X^T X)^{-1} X^T y \sim N(\beta, (X^T X)^{-1} \sigma^2)$$

smallest var among unbiased estimators

$$\mathbf{My}_m = k^T C_n^{-1} y = k^T (K + \sigma^2 I)^{-1} y$$

$$\sigma_{y_m}^2 = c - k^T C_n^{-1} k$$

$$\mathbb{E}[\text{Ex}(x)(\hat{f}(x) - Y)^2] = \mathbb{E}_\theta[(\hat{f}(x) - \mathbb{E}_\theta[\hat{f}(x)])^2]$$

$$+ (\mathbb{E}[\hat{f}(x) - \mathbb{E}[Y|X=x]])^2 + \mathbb{E}[(Y - \mathbb{E}[Y|X=x])^2]$$

$$p(a_1=2) = N(a_2 | \text{wt} \sum_i \tilde{z}_i^2 (z_i^2 - 1), 2 \sum_i -\tilde{z}_i \tilde{z}_i^T \tilde{z}_i^2).$$

$$\text{Ridge Regression}:$$

$$\text{Regularization}, CV, Ensembles$$

$$\text{LASSO}:$$

$$RSS(\beta; \lambda) = (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta$$

$$\text{Laplace prior: } p(\beta) = \frac{1}{4\sigma} \exp(-\|\beta\|_1)$$

$$\text{Gaussian prior: } \beta \sim N(0, \sigma^2 \lambda I).$$

$$\text{CV / Bootstrap / Jackknife: }$$

$$\text{Naive Bootstrap: } \hat{R}^{(1)} = \frac{1}{B} \sum_{b=1}^B (y_b - \hat{f}^{(b)}(x))^2$$

$$\text{Vb P.S. in } \mathbb{E}[\hat{f}^{(b)}] = 1 - (1 - \frac{1}{n})^n \rightarrow 1 - \frac{1}{e} = 0.632$$

$$\text{0.632 bootstrap: } \hat{R}_{\text{cv}} = 0.818 \text{ Ptrain + 0.18 } \hat{R}^{(1)}$$

$$\text{Jackknife: } 100 \text{ estimator } \hat{s}_{n-1} \rightarrow \text{est. bias}$$

$$\mathbb{E}[\hat{s}_n] - s = \frac{a_1}{n} + \frac{a_2}{n} + \dots$$

$$\mathbb{E}[\hat{s}_{n-1}] - s = \frac{a_1}{n-1} + \frac{a_2}{n-1} + \dots$$

$$\text{bias} \triangleq (n-1) (\hat{s}_n - \hat{s}_{n-1}) \hat{s}_0 = \frac{1}{n} \sum_{i=1}^n \hat{s}_{n-i}$$

$$\mathbb{E}[\hat{s}_{n-1}] = \frac{a_1}{n-1} + \frac{(n-1)\hat{s}_0}{(n-1)n_2 + 0(n^{-2})}$$

$$\text{JK estimator: } \hat{s}_{JK} = \hat{s}_n - \text{bias JK}$$

$$\text{st. } (w_2^T y_i + w_3^T o) - \max(w_2^T y_i)$$

$$\text{Structural Sums: } \sum_{i=1}^n \sum_{j=1}^n \delta_{ij} \delta_{ij}$$

$$\text{Multiclass SVMs: } \min_{w \geq 0} \frac{1}{2} w^T w + C \sum_{i=1}^n \delta_{i, y_i}$$

$$\text{hard fractional margin: } \frac{1}{2} w^T w + C \sum_{i=1}^n \delta_{i, y_i}$$

$$\text{margin: } \min_{w \geq 0} \frac{1}{2} w^T w + C \sum_{i=1}^n \delta_{i, y_i}$$

$$\text{Ensemble Methods}$$

$$\text{Prediction by Gaussian Process}$$

$$\mathbb{P}(\mathbf{y} | \mathcal{X}) = \mathbb{P}(\mathbf{y} | \mathbf{x}_1, \dots, \mathbf{x}_n)$$

② Feasible:  $g_i(w^*) = 0$ ; ③ (Feasible)  $\mathbb{P}(\mathbf{y} | \mathcal{X})$

④ Complementary Slackness:  $\alpha_i^* g_i(w^*)$

⑤ (Minimizes Lagrangian)  $\frac{\partial}{\partial w} L(w, \lambda^*, \alpha^*)$

Hard Margin SVM minimize ||w||

minimize  $T(w) = \frac{1}{2} w^T w$

subject to  $\forall i: y_i(w^T \mathbf{x}_i + w_0) \geq 1$

subject to  $y_i: \alpha_i \geq 0$ ,  $\frac{1}{n} \sum_{i=1}^n \alpha_i = 1$

$w^* = \frac{1}{n} \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i$

$\mathbb{P}(\mathbf{y} | \mathcal{X}) = \prod_{i=1}^n \mathbb{P}(y_i | \mathbf{x}_i)$

Optimal margin:  $w^T w = \frac{1}{2} w^T w + C \sum_{i=1}^n \delta_{i, y_i}$

subject to  $y_i: \delta_i \geq 0$

$\mathbb{P}(\mathbf{y} | \mathcal{X}) = \prod_{i=1}^n \mathbb{P}(y_i | \mathbf{x}_i)$

maximize  $W(w) = \frac{1}{2} w^T w + C \sum_{i=1}^n \delta_{i, y_i}$

subject to  $\forall i: y_i(w^T \mathbf{x}_i + w_0) \geq 1$

$\mathbb{P}(\mathbf{y} | \mathcal{X}) = \prod_{i=1}^n \mathbb{P}(y_i | \mathbf{x}_i)$

minimize  $T(w) = \frac{1}{2} w^T w$

subject to  $\forall i: y_i(w^T \mathbf{x}_i + w_0) \geq 1$

$\mathbb{P}(\mathbf{y} | \mathcal{X}) = \prod_{i=1}^n \mathbb{P}(y_i | \mathbf{x}_i)$

minimize  $\frac{1}{2} w^T w + C \sum_{i=1}^n \delta_{i, y_i}$

subject to  $\forall i: y_i(w^T \mathbf{x}_i + w_0) \geq 1$

$\mathbb{P}(\mathbf{y} | \mathcal{X}) = \prod_{i=1}^n \mathbb{P}(y_i | \mathbf{x}_i)$

minimize  $\frac{1}{2} w^T w + C \sum_{i=1}^n \delta_{i, y_i}$

subject to  $\forall i: y_i(w^T \mathbf{x}_i + w_0) \geq 1$

$\mathbb{P}(\mathbf{y} | \mathcal{X}) = \prod_{i=1}^n \mathbb{P}(y_i | \mathbf{x}_i)$

minimize  $\frac{1}{2} w^T w + C \sum_{i=1}^n \delta_{i, y_i}$

subject to  $\forall i: y_i(w^T \mathbf{x}_i + w_0) \geq 1$

$\mathbb{P}(\mathbf{y} | \mathcal{X}) = \prod_{i=1}^n \mathbb{P}(y_i | \mathbf{x}_i)$

minimize  $\frac{1}{2} w^T w + C \sum_{i=1}^n \delta_{i, y_i}$

subject to  $\forall i: y_i(w^T \mathbf{x}_i + w_0) \geq 1$

$\mathbb{P}(\mathbf{y} | \mathcal{X}) = \prod_{i=1}^n \mathbb{P}(y_i | \mathbf{x}_i)$

Elastic Net SVM

KKT Conditions

**Bagging:** for  $b=1$  to  $B$  do

$\hat{z}^{*b} = b\text{-th bootstrap sample from } \mathcal{Z}$

Construct classifier  $c_b$  based on  $\hat{z}^{*b}$ .

end for

return  $\hat{c}_B(x) = \text{sign}(\sum_{i=1}^B c_i(x))$ .

Covariance  $\mathbb{V}$ , Variance similar. Biases -

**AdaBoost**: initialize weight  $w_i = \frac{1}{n}$ .

for  $b=1$  to  $B$  do

Fit  $c_b(x)$  to the training data with  $w_i$

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

$\alpha_b \leftarrow \log \frac{1-\hat{e}_b}{\hat{e}_b}$

Set  $w_i' : w_i \leftarrow w_i \exp(\alpha_b \mathbb{I}\{y_i \neq c_b(x_i)\})$

end for.

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

**PAC Learning**

Learnable:  $\exists A$  can learn any concept in  $\mathcal{C}$ .

$\forall A$  receives sample  $n \geq \text{poly}(\lambda, \epsilon, \delta)$ ,

$\text{dim}(\mathcal{X}) \cdot P(C(R)) - \text{int}(R_C) \leq \epsilon \geq 1 - \delta$

Efficient PAC learnable:  $A$  runs in time

$\text{poly}(\lambda, \epsilon, 1/\delta)$ . unions of  $k$  intervals  $\leq 2k$

**VC dimension**: convex polygons in  $\mathbb{R}^2$

with at most  $K$  vertices:  $2K+1$

**VC inequality**:  $\Pr[\hat{C}(R) \geq \epsilon] \leq \Pr[\text{sup}_{f \in \mathcal{F}} |f(x)| - R(f) > \frac{\epsilon}{2}]$

**Hoeffding's Lemma**:  $\mathbb{E}[\text{sup}_{x \in \mathcal{X}} f(x)] \leq \exp((\epsilon^2/2)(n-1)/2) \leq \exp(-\epsilon^2/2)$

**Markov's Ineq.**:  $P\{X \geq t\} \leq \frac{\mathbb{E}[X]}{t}$

**Hoeffding's Lemma**:  $\mathbb{E}[\text{sup}_{x \in \mathcal{X}} f(x)] \leq \exp((\epsilon^2/2)(n-1)/2) \leq \exp(-\epsilon^2/2)$

**Uniform Bound**:  $\Pr\{\cup_{i=1}^n A_i\} = \sum_{i=1}^n \Pr\{A_i\}$

**Hoeffding's Bound**:  $\Pr\{\cup_{i=1}^n A_i\} > \epsilon \leq 2^n \exp(-\epsilon^2/2)$

**Nonparametric Bayesian Methods**

Dirichlet dist.  $D_{\alpha}(\mathbf{x}|\alpha) = \frac{1}{B(\alpha)} \prod_{i=1}^B X_i^{\alpha_i-1}$

**AdaBoost**:  $\Pr[x_i \in \mathcal{C}] = \frac{\Pr[x_i \in \mathcal{C}]}{\Pr[x_i \in \mathcal{C} \cup \mathcal{C}^c]} = \frac{\Pr[x_i \in \mathcal{C}]}{1 + \Pr[x_i \in \mathcal{C}^c]}$

**Naive Bayes**:  $\Pr[x_i \mid \mathcal{C}]$

**DP Mixture Models**

**EM**: likelihood  $\Pr(X|\theta) = \prod_{i=1}^n \Pr(x_i|\theta_i)$  ② **Stepwise reg.**:  $F_{\text{new}} \leftarrow F_{\text{old}} + \alpha \mathbb{V}$

① Prob. of clusters:  $p = (p_1, p_2, \dots) \sim \text{Geom}(\mu)$ .

② Centers of clusters:  $\mu_k \sim N(\mu_0, \sigma_0)$ ,  $k=1, 2, \dots$

③ Assignment:  $z_{ik} \sim \text{Categorical}(p_i), i=1, \dots, n$

○ Cond. of data:  $x_{ik} \sim N(\mu_{z_{ik}}, \sigma^2), i=1, \dots, n$

**EM**: sampling  $\Pr(z_{ik}=k \mid z_{-ik}, x_{ik}, \mu_i, \sigma^2) \propto$

$= \frac{\Pr(z_{ik}=k) \Pr(x_{ik} \mid z_{ik}=k, \mu_i, \sigma^2)}{\Pr(z_{ik}=1) \Pr(x_{ik} \mid z_{ik}=1, \mu_1, \sigma^2) + \dots + \Pr(z_{ik}=K) \Pr(x_{ik} \mid z_{ik}=K, \mu_K, \sigma^2)}$

$= \begin{cases} \frac{p_{ik}}{\sum_{i=1}^K p_{ik}} & \Pr(x_{ik} \mid z_{ik}=k, \mu_k, \sigma^2) \\ 0 & \text{otherwise} \end{cases}$

**EM**:  $\Pr(z_{ik}=k \mid z_{-ik}, x_{ik}, \mu_i, \sigma^2) \propto$

$= \begin{cases} \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x_{ik}-\mu_i)^2) & \Pr(x_{ik} \mid z_{ik}=k, \mu_k, \sigma^2) \\ 0 & \text{otherwise} \end{cases}$

**EM**:  $\Pr(z_{ik}=k \mid z_{-ik}, x_{ik}, \mu_i, \sigma^2) \propto$

$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x_{ik}-\mu_i)^2) \Pr(x_{ik} \mid z_{ik}=k, \mu_k, \sigma^2)$

**EM**:  $\Pr(z_{ik}=k \mid z_{-ik}, x_{ik}, \mu_i, \sigma^2) \propto$

$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x_{ik}-\mu_i)^2) \Pr(x_{ik} \mid z_{ik}=k, \mu_k, \sigma^2)$

**EM**:  $\Pr(z_{ik}=k \mid z_{-ik}, x_{ik}, \mu_i, \sigma^2) \propto$

$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x_{ik}-\mu_i)^2) \Pr(x_{ik} \mid z_{ik}=k, \mu_k, \sigma^2)$

**EM**:  $\Pr(z_{ik}=k \mid z_{-ik}, x_{ik}, \mu_i, \sigma^2) \propto$

$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x_{ik}-\mu_i)^2) \Pr(x_{ik} \mid z_{ik}=k, \mu_k, \sigma^2)$

**EM**:  $\Pr(z_{ik}=k \mid z_{-ik}, x_{ik}, \mu_i, \sigma^2) \propto$

$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x_{ik}-\mu_i)^2) \Pr(x_{ik} \mid z_{ik}=k, \mu_k, \sigma^2)$

**EM**:  $\Pr(z_{ik}=k \mid z_{-ik}, x_{ik}, \mu_i, \sigma^2) \propto$

$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x_{ik}-\mu_i)^2) \Pr(x_{ik} \mid z_{ik}=k, \mu_k, \sigma^2)$

**EM**:  $\Pr(z_{ik}=k \mid z_{-ik}, x_{ik}, \mu_i, \sigma^2) \propto$

$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x_{ik}-\mu_i)^2) \Pr(x_{ik} \mid z_{ik}=k, \mu_k, \sigma^2)$

**EM**:  $\Pr(z_{ik}=k \mid z_{-ik}, x_{ik}, \mu_i, \sigma^2) \propto$

$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x_{ik}-\mu_i)^2) \Pr(x_{ik} \mid z_{ik}=k, \mu_k, \sigma^2)$

**EM**:  $\Pr(z_{ik}=k \mid z_{-ik}, x_{ik}, \mu_i, \sigma^2) \propto$

$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x_{ik}-\mu_i)^2) \Pr(x_{ik} \mid z_{ik}=k, \mu_k, \sigma^2)$

**EM**:  $\Pr(z_{ik}=k \mid z_{-ik}, x_{ik}, \mu_i, \sigma^2) \propto$

$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x_{ik}-\mu_i)^2) \Pr(x_{ik} \mid z_{ik}=k, \mu_k, \sigma^2)$

**EM**:  $\Pr(z_{ik}=k \mid z_{-ik}, x_{ik}, \mu_i, \sigma^2) \propto$

$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp(-\frac{1}{2\sigma^2}(x_{ik}-\mu_i)^2) \Pr(x_{ik} \mid z_{ik}=k, \mu_k, \sigma^2)$

**EM**: likelihood  $\Pr(X|\theta) = \prod_{i=1}^n \Pr(x_i|\theta_i)$  ② **Stepwise reg.**:  $F_{\text{new}} \leftarrow F_{\text{old}} + \alpha \mathbb{V}$

$\Pr(X|\theta) = \prod_{i=1}^n \Pr(x_i|\theta_i)$

$E\text{-step}$ :  $\gamma_{ik} = \Pr(x_i \mid z_{ik}, \theta_i)$

$\Pr(X|\theta) = \prod_{i=1}^n \Pr(x_i \mid z_{ik}, \theta_i)$

**EM**: likelihood  $\Pr(X|\theta) = \prod_{i=1}^n \Pr(x_i|\theta_i)$  ② **Stepwise reg.**:  $F_{\text{new}} \leftarrow F_{\text{old}} + \alpha \mathbb{V}$

$\Pr(X|\theta) = \prod_{i=1}^n \Pr(x_i|\theta_i)$

$J(F_{\text{new}}) = \mathbb{E}_{\Pr[X|F_{\text{new}}]} [\exp(-\gamma F_{\text{new}}) - \alpha \mathbb{V}]$

$= \mathbb{E}_{\Pr[X|F_{\text{new}}]} [\exp(-\gamma F_{\text{new}}) - \alpha \mathbb{V}]$