

Linear Regression

$$f(x) = w_0 + w_1 x = w_0 + \sum_{i=1}^n w_i x_{i,1} = w_0 + w^T x$$

$$L(f) = \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i) \quad (\rightarrow \text{MSE w/ } l_{\text{square}})$$

$$\text{Loss functions:}$$

- $l_{\text{square}} = (f(x) - y)^2$
- $l_{\text{huber}} = \begin{cases} \frac{1}{2} (y - f(x))^2 & \text{for } |y - f(x)| \leq \delta \\ \delta (|y - f(x)| - \frac{1}{2} \delta) & \text{else} \end{cases}$

Closed Form Solution: $\nabla_w L(\hat{w}) = 0$

$$\text{Normal Equations: } (X^T X)^{-1} X^T y$$

where X is $n \times (d+1)$, $X = \begin{bmatrix} 1 & x_1 & \dots & x_d \end{bmatrix}$

$$\text{Gradient Descent: } w^{t+1} = w^t - \eta \nabla_w L(w^t)$$

- e.g. stop when $|L(w^{t+1}) - L(w^t)| \leq \epsilon$
- momentum: $w^{t+1} - w^t = \alpha(w^t - w^{t-1}) - \eta \nabla_w L(w^t)$

Convexity, order conditions:

$$0 \ L(\lambda w + (1-\lambda)v) \leq \lambda L(w) + (1-\lambda)L(v)$$

$$1 \ L(v) \geq L(w) + \nabla L(w)^T (v-w)$$

$$2 \text{ Hessian } \nabla^2 L(w) \geq 0$$

Convexity preserving operations

- $\alpha f + \beta g$ for $\alpha, \beta \geq 0$ and f, g convex
- $f \circ g = f(g(x))$ if f convex, g affine
- $\max \{f(x), g(x)\}$ if f, g convex

Bias: dist. of avg. model to ground truth

$$\mathbb{E}_x \left[(f^*(x) - \frac{1}{n} \sum_{j=1}^n \hat{f}_j(x))^2 \right]$$

Variance: avg. dist. of models to ground truth

$$\mathbb{E}_x \left[\frac{1}{n} \sum_{j=1}^n (\hat{f}_j(x) - \bar{f}(x))^2 \right]$$

Regularization: focus on low degrees

$$\text{Lasso: } \arg \min_w \|y - \Phi w\|_2^2 + \lambda \|w\|_1$$

$$\text{Ridge: } \arg \min_w \|y - \Phi w\|_2^2 + \lambda \|w\|_2^2 \quad \text{all degrees}$$

Matrix is psd: All EVs ≥ 0

- All principal minors are non-neg.
- remove rows/cols w/ same index

For pd only leading minors > 0

$$\log(u \cdot v) = \log(u) + \log(v) / \alpha^b = e^{b \log(a)}$$

$$\log(u^v) = v \cdot \log(u) / \log_b(u) = \log_b(u)^v$$

Classification

Linear Classifier: $f(x) = w^T x$ for some $w \in \mathbb{R}^d$ Non-linear: $f(x) = w^T \phi(x)$, $\phi \in \mathbb{R}^d \rightarrow \mathbb{R}^p$

Loss functions:

$g_{\text{exp}}(y, \hat{f}(x)) = e^{-\hat{f}(x)}$

$g_{\text{log}}(y, \hat{f}(x)) = \log(1 + e^{-\hat{f}(x)})$

$g_{\text{lin}}(y, \hat{f}(x)) = -y \hat{f}(x)$

less sensitive to outliers

$y \hat{f}(x)$

$\text{Cross-Entropy Loss}$

$l_{\text{ce}}(f(x), y) = -\log \left(\frac{e^{f(y)}}{\sum_i e^{f(y_i)}} \right)$

$\hat{f}(x) = f(x) \text{ where } y = \text{margin}(w)$

$$\text{Hard-Margin SVM } w_{\text{HM}} = \arg \max_{\|w\|_2} \min_i y_i \langle w, x_i \rangle$$

$$\text{w}_{\text{HM}} \parallel \arg \min_{\|w\|_2} \text{st. } y_i w^T x_i \geq 1 \quad \forall i = 1, \dots, n$$

Null Hypothesis: the critical class, want no error

Type 1 Error: False Pos / Type 2 Error: FN

$$\text{FNR} = \frac{\# \text{FN}}{\# P} / \text{FPR} = \frac{\# \text{FP}}{\# N} / \text{TPR} = \frac{\# \text{TP}}{\# P} = 1 - \text{FNR}$$

Asymm. lin. classifier: $C_{\text{FN}} \cdot \text{FNR} + C_{\text{FP}} \cdot \text{FPR}$ ROC curve: $\text{TPR} \uparrow \quad \text{FPR} \downarrow$

we want top left

$$\text{Precision: } \frac{\# \text{TP}}{\# [\hat{y}=1]} \approx \mathbb{P}_n [\hat{y}=1 \mid \hat{y}=1]$$

$$\text{Recall: } \frac{\# \text{TP}}{\# [\hat{y}=1]} \approx \mathbb{P}_n [\hat{y}=1 \mid y=1]$$

$$\text{False Discovery Rate: } \frac{\# \text{FP}}{\# [\hat{y}=1]} \approx \mathbb{P}_n [y=1 \mid \hat{y}=1]$$

$$\text{FPR: } \frac{\# \text{FP}}{\# [\hat{y}=-1]} \approx \mathbb{P}_n [\hat{y}=1 \mid y=-1]$$

$$\text{F1 score: } \frac{2}{\frac{1}{\text{recall}} + \frac{1}{\text{precision}}} \rightarrow \text{forces both to be large}$$

Asymmetric loss: False Positives False Neg.

$$l(\hat{y}(x), y) = C_{\text{FP}} \mathbb{I}_{\hat{y}(x)=1, y=-1} + C_{\text{FN}} \mathbb{I}_{\hat{y}(x)=-1, y=1}$$

$$(AB \dots)^{-1} = \dots B^{-1} A^{-1}$$

$$(AB \dots)^T = \dots B^T A^T$$

$$\text{Tr}(A) = \sum_i \lambda_i$$

$$\text{Tr}(AB) = \text{Tr}(BA)$$

$$\text{Tr}(A+B) = \text{Tr}(A) + \text{Tr}(B)$$

$$a^T a = \text{Tr}(a a^T)$$

$$\det(A) = \prod_i \lambda_i$$

$$\det(cA) = c^n \cdot \det(A)$$

$$\det(AB) = \det A \cdot \det B$$

$$g(f(x)) = J_f(f(x)) \cdot f'(x)$$

$$= \frac{\partial g}{\partial f} \cdot \frac{\partial f}{\partial x}$$

Kernels

$$\text{Solve } \hat{w} = \arg \min_{w \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n l(y_i, w^T \phi(x_i))$$

$$\text{Output } \hat{f}(x) = \sum_i \hat{w}_i \phi(x_i) = \alpha \Phi \phi(x)$$

$$1) \text{Solve } \hat{\alpha} = \arg \min_{\alpha \in \mathbb{R}^n} \frac{1}{n} \sum_{i=1}^n l(y_i, \alpha^T \Phi \phi(x_i))$$

$$2) \hat{f}(x) \text{ only depends on } x \text{ via inner products}$$

$$3) \text{Replace } \langle \cdot, \cdot \rangle \text{ w/ } k(x, z) \rightarrow \hat{f}(x) = \sum_i \hat{\alpha}_i k(x, x_i)$$

$$\Phi = \begin{pmatrix} -\phi(x_1) \\ \vdots \\ -\phi(x_n) \end{pmatrix} \quad / \quad \Phi \cdot \phi(x_i) = \begin{pmatrix} \langle \phi(x_1), \phi(x_i) \rangle \\ \vdots \\ \langle \phi(x_n), \phi(x_i) \rangle \end{pmatrix}$$

$$\text{Kernel function } k: X \times X \rightarrow \mathbb{R}$$

$$\text{Kernel matrix } K = \begin{pmatrix} k(x_1, x_1) & \dots & k(x_1, x_n) \\ \vdots & \ddots & \vdots \\ k(x_n, x_1) & \dots & k(x_n, x_n) \end{pmatrix}$$

$$\text{Kernelized regression: } \hat{\alpha} = \arg \min_{\alpha \in \mathbb{R}^n} \|y - K\alpha\|^2$$

$$\text{w/ } \hat{f}(x) = \hat{\alpha}^T K x \Rightarrow \hat{\alpha} = K^{-1} y$$

Valid kernels: symmetric $k(x, z) = k(z, x)$ • k is pos. semidefinite (all $\lambda \geq 0$)⇒ e.g. $h(\langle x, z \rangle)$, $h(\|x - z\|)$ Kernel composition: if k_1, k_2 valid kernels

$$\cdot k(\langle y, v \rangle) = k_1(x, u) [+, \cdot] k_2(y, v)$$

$$\cdot k(x, y) = k_1(x, y) [+, \cdot] k_2(x, y)$$

Dimension Reduction compressed

$$\text{PCA: } (w^*, z^*) = \arg \min_{(w, z)} \sum_{i=1}^n \|z_i^T w - x_i\|^2$$

$$\Rightarrow w^* = \arg \min_{\|w\|_2=1} \sum_{i=1}^n \|w^T x_i - x_i\|^2 \quad (z_i^* = w^T x_i)$$

$$\stackrel{M=0}{=} \arg \max_{\|w\|_2=1} w^T \sum_i w_i \quad \sum_i w_i = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

⇒ solution is principal EV of Σ

$$\text{If } \sum_i \lambda_i v_i v_i^T, \quad w = (v_1 \dots v_k), \quad z_i = w^T x_i$$

↳ if $X \in \mathbb{R}^{n \times d} = USV^T$, first k cols of V are sol

Kernel PCA:

$$\hat{\alpha} = \arg \max_{\alpha^T K \alpha = 1} \alpha^T K^T K \alpha = \arg \max_{\alpha^T K \alpha} \frac{\alpha^T K^T K \alpha}{\alpha^T \alpha}$$

$$\text{Let } K = \sum_{i=1}^n \lambda_i v_i v_i^T. \quad \text{Then } \alpha^{(i)} = \frac{1}{\sqrt{\lambda_i}} v_i$$

$$\text{and } z_i = \sum_{j=1}^n \alpha_j^{(i)} k(x_j, x_i)$$

$$\text{Autoencoders: } \min_w \sum_{i=1}^n \|x_i - f(x_i; w)\|^2$$

$$\frac{1}{n} \sum_{i=1}^n \|x_i - f(x_i; w)\|^2$$

$$= \frac{1}{n} \sum_{i=1}^n \|x_i - \sum_j w_{ij} \phi_j(x_i)\|^2$$

$$= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^p w_{ij}^2 \phi_j(x_i)^2$$

$$= \frac{1}{n} \sum_{j=1}^p w_{j,1}^2 \sum_{i=1}^n \phi_j(x_i)^2$$

$$= \frac{1}{n} \sum_{j=1}^p w_{j,1}^2 \text{Tr}(\Phi \Phi^T)$$

Clustering

k-means algorithm $O(n \cdot k \cdot d)$

- 1) init $\mu^{(0)} = [M_1^{(0)}, \dots, M_k^{(0)}]$
- 2) until converge: $z_i \leftarrow \operatorname{arg\min}_j \|x_i - M_j^{(t-1)}\|_2$
 $M_j^{(t)} \leftarrow \frac{1}{n_j} \sum_{i:z_i=j} x_i$ (\rightarrow mean)

k-means++: $M_1^{(0)}$ = random x_i
 add other centers proportional to squared dist. to closest selected
 \hookrightarrow cost $O(\log k)$ times optimal

Gaussian Mixture Models

Hard-EM: 1) init parameters $\theta^{(0)}$
 2) E-step: $z_i^{(t)} = \operatorname{arg\max}_z P(z|x_i, \theta^{(t-1)})$
 $= \operatorname{arg\max}_z P(z|\theta^{(t-1)}) \cdot P(x_i|z, \theta^{(t-1)})$
 3) $D^{(t)} = \{(x_1, z_1^{(t)}), \dots, (x_n, z_n^{(t)})\}$
 4) M-step: Compute MLE as for GBC
 $\theta^{(t)} = \operatorname{arg\max}_{\theta} P(D^{(t)} | \theta)$

Soft-EM: 1) init $\mu^{(0)}, \Sigma^{(0)}, w^{(0)}$
 2) E-step: calc cluster probabilities
 $\gamma_j^{(t)}(x_i)$ given estimates of $(t-1)$
 3) M-step: Fit clusters to data points
 $w_j^{(t)} \leftarrow \frac{1}{n} \sum_{i=1}^n \gamma_j^{(t)}(x_i)$
 $\gamma_j^{(t)} \leftarrow \frac{\sum_{i=1}^n \gamma_j^{(t)}(x_i) \cdot x_i}{\sum_{i=1}^n \gamma_j^{(t)}(x_i)}$
 $\Sigma_j^{(t)} \leftarrow \frac{\sum_{i=1}^n \gamma_j^{(t)}(x_i)(x_i - \mu_j^{(t)})(x_i - \mu_j^{(t)})^T}{\sum_{i=1}^n \gamma_j^{(t)}(x_i)}$

Constrained GMMs: # params
 Spherical $\Sigma_i = \sigma_i^2 \cdot I$ k

Diagonal $\Sigma_i = \operatorname{diag}(\sigma_{i,1}^2, \dots, \sigma_{i,d}^2)$ $k \cdot d$

Tied $\Sigma_1 = \dots = \Sigma_k \rightarrow$ All shapes the same $\frac{d(d+1)}{2}$

Full w, μ, Σ arbitrary $k \cdot \frac{d(d+1)}{2}$

Initialization: • weights: uniform dist.

• means: random or k-means++

• variances: spherical

Semi-supervised: $\gamma_j^{(t)}(x_i) = [j = y_i]$
 for points w/ label. Rest Soft-EM

Probabilistic Modelling

Gaussian MLE, $p(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
 log-likelihood $LLH(\mu, \sigma) = \log p(D; \mu, \sigma)$
 $\nabla_{\mu} LLH = 0 \Rightarrow \hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$
 $\nabla_{\sigma} LLH = 0 \Rightarrow \hat{\sigma}_{MLE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_{MLE})^2}$
 $\Rightarrow P = N(\hat{\mu}_{MLE}, \hat{\sigma}_{MLE})$

Discriminative: estimate $p(y|x)$

Decision boundaries btw. classes

Generative: estimate $p(y, x) = p(x|y) \cdot p(y)$

Model prob. distribution of data

Gaussian Naive Bayes: Decision boundary

linear, if Σ_y diagonal and same for all y

Finding MLE for $P_{x|y}$

$\hat{\mu}_y = \frac{1}{\#y=y} \sum_{i:y_i=y} x_i, \hat{\sigma}_{y,k} = \sqrt{\frac{1}{\#y=y} \sum_{i:y_i=y} (x_{i,k} - \hat{\mu}_{y,k})^2}$

Bayesian Modeling: We know something about

the distribution of $\theta \rightarrow p(D) = \int p(D|\theta) \cdot p(\theta) d\theta$

\Rightarrow output $\operatorname{arg\max}_{\theta} p(D, \theta) = p(D|\theta) \cdot p(\theta)$

Generative Modeling for classification

1) estimate prior on labels $p(y)$

2) estimate conditional dist. $p(x|y)$

3) using Bayes rule: $p(y|x) = \frac{1}{Z} p(y) \cdot p(x|y)$
 with $Z = \sum_y p(y) \cdot p(x|y)$

MLE for Gaussian Bayes Classifier

• Given dataset $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$

• $P(Y=y) = \hat{p}_y = \frac{\text{Count}(Y=y)}{n}$

• $P(x; \hat{\mu}_y, \hat{\Sigma}_y) \Rightarrow \hat{\mu}_y = \frac{1}{\text{Count}(Y=y)} \sum_{i:y_i=y} x_i$
 $\Rightarrow \hat{\Sigma}_y = \frac{1}{\text{Count}(Y=y)} \sum_{i:y_i=y} (x_i - \hat{\mu}_y)(x_i - \hat{\mu}_y)^T$

$P(X, Y) = P(Y) \cdot P(X|Y)$

Transformers meaning predicted

Word $i \quad v_i \quad \downarrow \quad q_i \quad \downarrow \quad k_i$

music [0, 1, 0] [0, 0] [1, 1]

rock [2, 0, 6] [1, 0] [0, 1]

like [0, 0, 0] [0, 1] [0, 0]

score_{i,j} $\propto q_i \cdot k_j^T$ attending attended to
 \hookrightarrow want score_{i,j} ≈ 0 if $j > i$