

REGRESSION

Linear

$$L(w) = \frac{1}{n} \sum_i (y_i - w^T x_i)^2 = \frac{1}{n} \|y - w x\|_2^2$$

$$\nabla_w L(w) = -\frac{2}{n} (w^T x_i - y_i) x_i = -\frac{2}{n} (x^T x w - x^T y) \in \mathbb{R}^d$$

$$\nabla_w^2 L(w) = \frac{2}{n} x^T x$$

$$\hat{w} = (x^T x)^{-1} x^T y$$

$\text{det}(A^T A) > 0$ $\Rightarrow A^T A$ is invertible

Lasso (convex, + strongly if X is full rank)

$$L(w) = \|y - \Phi w\|_2^2 + \lambda \|w\|_1$$

Ridge (strictly convex, even if X not full rank)

$$L(w) = \|y - \Phi w\|_2^2 + \lambda \|w\|_2^2$$

$\hat{w} = (X^T X + \lambda I_d)^{-1} X^T y$

SD directly better conditioned as changing λ_{\max} and λ_{\min}

OPTIMIZATION

Gradient Descent (conv. to stationary pt)

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

Learning rate

$$\|w^t - w^*\|_2 \leq \rho^t \|w^0 - \hat{w}\|_2$$

$$\rho = \|I - \eta X^T X\|_{\text{op}}$$

$$\eta < \frac{\lambda_{\min}(X^T X)}{\lambda_{\max}(X^T X)}$$

$$\eta_{\text{opt}} = \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min}} \Rightarrow \rho_{\text{opt}} = \frac{\lambda_{\max} - \lambda_{\min}}{\lambda_{\max} + \lambda_{\min} + 1}$$

Momentum: dampens oscillations

$$w^{t+1} = w^t + \alpha (w^t - w^{t-1}) - \eta \nabla_w L(w^t)$$

Adaptive: w_{it} changed + lot \rightarrow smaller stepsize

$$w^{t+1} = w^t - \eta_t \nabla_w L(w^t)$$

Training loss might increase at some iter.

Despite in t sampled uniformly w replacement

SGD update is a descent direction in expectation

Large batch size \rightarrow small variance

Escape saddle points

Fewer steps than GD, but may take more and also not converge to right sol, but swirl

Converge? (loc \Rightarrow glb)

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

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

$$J^2(w) \geq 0$$

Strong convexity \rightarrow uniqueness

Complexity

GD: Steps \propto req. to satif. $\|w^t - w^*\| < \epsilon$

$$\gamma = \frac{\log \epsilon}{\log \frac{1}{\rho}}$$

Total: $\sigma(n^2 + d^2)$; Single: $\sigma(nd)$

CF: $\sigma(d^3)$; lin. sys. of eqs: $\sigma(n^2)$

MODEL SELECTION

Expected estimation error

$$E_x[\ell(f_\theta(x), f^*(x))]$$

Generalization error

$$L(f; P_{\mathcal{D}_x}) = E_{\mathcal{D}_x}[\ell(f_\theta(x), y)]$$

$$= \sum_y \sum_x \ell(f_\theta(x), y) P(\bar{x}, y)$$

$$= \int \ell(f_\theta(x), y) p(x, y) dx dy$$

$$= E_x[L(\ell(f_\theta(x), f^*(x)))] + \sigma^2$$

Training error

$$\frac{1}{|\mathcal{D}_x|} \sum_{(x,y) \in \mathcal{D}_x} \ell(f_\theta(x), y)$$

Cross validation

For all folds $K=1, \dots, K$:

Train model $f_{\theta, k}$ on D_{k-1}

Validation error $L_k^* = \frac{1}{|\mathcal{D}_k|} \sum_{(x,y) \in \mathcal{D}_k} \ell(f_{\theta, k}(x), y)$

Cross validation error $CV_k(M) = \frac{1}{K} \sum_k L_k(M)$

Model selection: Pick M^* with lowest $CV_k(M)$

Model training: Train model $\hat{f} = f_{\theta, M^*}$, D_{k-1}

Model evaluation: Estimate gen. error using D_{k-1}

LOOCV: $K=|\mathcal{D}_{k-1}|$ (deterministic)

$\cdot M(D_{k-1}) \approx M(D_{k-1})$

\cdot Poor L_k^* , but $CV_k(M)$ may be good gen. est.

\cdot Intensive computationally (M Daniel models)

BIAS-VARIANCE TRADEOFF

Bias

$$\text{Bias}(\hat{f}_0) = E_x[\ell(\hat{f}_0(x) - f^*(x))^2]$$

$$= \frac{1}{n} \sum_{i=1}^n \hat{f}_0(x_i) - f^*(x_i)$$

$$\in \mathbb{R}$$

$\hat{f}_0(x) = \hat{w}^T x$

\hat{w} is $(X^T X)^{-1} X^T y$

\hat{w} is invertible

\hat{w} is $(X^T X)^{-1} X^T y$

\hat{w} is invertible

MAXIMUM LIKELIHOOD ESTIM.

Likelihood of the data for given distr. $P(\cdot; \theta)$

$$P(D; \theta) = \prod_{i=1}^n P(x_i, y_i; \theta) : P_x$$

$$= \prod_{i=1}^n p(y_i | x_i; \theta) : P_{y|x}$$

Maximum likelihood param. est.

$$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \log P(D; \theta)$$

Unsupervised learning, P_x

Gaussian distribution, $N(\mu, \sigma^2)$

$$\hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_{MLE})^2$$

Supervised learning, $P_{y|x}$

Gaussian noise model (regression)

Assume $y = w^T x + \epsilon$; $\epsilon \sim N(0, \sigma^2_\epsilon)$

$$\Rightarrow p(y|x; w) \sim N(w^T x, \sigma^2_\epsilon)$$

$$\text{squared loss} \quad \hat{w}_{MLE} = \underset{w}{\operatorname{argmin}} (y_i - w^T x_i)^2$$

Logistic noise model (classification)

$$\Rightarrow p(y|x; \theta) = \sigma(y \theta^T x) = \frac{1}{1 + e^{-y \theta^T x}}$$

$$\Rightarrow \hat{y}(x) = \underset{y \in \{-1, 1\}}{\operatorname{argmax}} \hat{p}(y|x)$$

$$\text{logistic loss} \quad \hat{\theta}_{MLE} = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmax}} \sum_{i=1}^n \log (1 - e^{-y_i \theta^T x_i})$$

MAXIMUM A POSTERIORI ESTIM.

Posterior distribution of θ given data

$$p(\theta|D) = \frac{p(D|\theta) \cdot p(\theta)}{p(D)}$$

$$P_x = \frac{\prod_{i=1}^n p(x_i; \theta) \cdot p(\theta)}{\int \prod_{i=1}^n p(x_i; \theta) \cdot p(\theta) d\theta} = \frac{\prod_{i=1}^n p(y_i | x_i, \theta) \cdot p(\theta)}{\int \prod_{i=1}^n p(y_i | x_i, \theta) \cdot p(\theta) d\theta}$$

Maximum a posteriori param. est.

$$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} \log p(\theta|D) = \underset{\theta}{\operatorname{argmax}} \log P(D; \theta)$$

Supervised distribution fun. P for P_{xy}

+ Gaussian prior: $w \sim \mathcal{N}(0, \sigma_w^2 I_d)$

$$\Rightarrow \hat{w}_{MAP} = \underset{w}{\operatorname{argmin}} \frac{1}{2\sigma_w^2} \|y - Xw\|_2^2 + \frac{1}{2\sigma_w^2} \|w\|_2^2$$

$$\lambda = \frac{\sigma_w^2}{\sigma_w^2} = \underset{w}{\operatorname{argmin}} \|y - Xw\|_2^2 + \frac{\sigma_w^2}{\sigma_w^2} \|w\|_2^2$$

+ Laplacian prior: $w \sim \text{Laplace}(0, \sigma_w)$

$$\Rightarrow \hat{w}_{MAP} = \underset{w}{\operatorname{argmin}} \frac{1}{2\sigma_w} \|y - Xw\|_2^2 + \frac{1}{\sigma_w} \|w\|_1$$

$$\lambda = \frac{2\sigma_w}{\sigma_w} = \underset{w}{\operatorname{argmin}} \|y - Xw\|_2^2 + \frac{2\sigma_w}{\sigma_w} \|w\|_1$$

Bayesian modeling averaging

Bayesian modeling allow outputting a distr. of possible θ

\hookrightarrow Instead of point est., posterior $p(\theta|D)$ automatically gives a distribution over distributions

Bayesian model averaging for P_x and $P_{y|x}$, respectively

$$\hat{p}(x|D) = \hat{E}_{\theta|D} [p(x|D)] = \int_{\theta} p(x|\theta) \hat{p}(\theta|D) d\theta$$

by choice of model computed from data

$$\hat{p}(y|x, D) = \hat{E}_{\theta|D} [p(y|x, \theta)] = \int_{\theta} p(y|x, \theta) \hat{p}(\theta|D) d\theta$$

Conjugate prior: yield a posterior within the same distr. as given likelihood function

GENERATIVE MODELING

Model $p(y, x) = p(x|y) \cdot p(y)$

Estimate prior on labels: $p(y)$

Estimate cond. distr. for each class y : $p(x|y)$

Pred. distr. using Bayes: $p(y|x) = \frac{p(y)}{p(x)} p(x|y)$

Decision rule: $y = \underset{y'}{\operatorname{argmax}} p(y'|x) = \underset{y'}{\operatorname{argmax}} p(y') p(x|y')$

$$= \underset{y'}{\operatorname{argmax}} \log p(y') + \log p(x|y')$$

Gaussian Naive Bayes (classification)

Naive: for each class label, each feature is indep.

$$\hookrightarrow \Sigma_y = \text{diag}(\Sigma_{y,1}, \dots, \Sigma_{y,d})$$

\hookrightarrow fail to capture correlation \rightarrow overly confident

Model class label as generated from cat. variable: $P(Y=y) = p_y$

Model features by (conditionally) independent Gaussians

$$p(x|y) = N(x; \mu_y, \Sigma_y)$$

MLE: $\hat{p}_{MLE} = \frac{\#\{y \neq y\}}{n}$

$$\hat{\mu}_{MLE} = \frac{1}{\#\{y \neq y\}} \sum_{i:y_i \neq y} x_i$$

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

Compare with log reg.: smaller uncertainty band, right shape

Gaussian Bayes classifier / QDA

Model class label as generated from cat. variable: $P(Y=y) = p_y$

Model features by multivariate Gaussians

$$p(x|y) = N(x; \mu_y, \Sigma_y)$$

MLE: $\hat{p}_y = \text{Count}(Y=y)$

$$\hat{\mu}_y = \frac{n}{\text{Count}(Y=y)} \sum_{i:y_i=y} x_i$$

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

Linear Discriminant analysis

$$\Sigma_y = \sigma^2 (\Sigma_1 = \dots = \Sigma_k) \text{ for all } y \Rightarrow \text{linear dbl}$$

Poisson Naive Bayes

$$\hat{p}_y = \frac{\text{Count}(Y=y)}{n}, \quad \lambda_{y,j} = \frac{\sum_{i:y_i=y} x_j^{(i)}}{\text{Count}(Y=y)}$$

Constrained GMMs (spherical, diff. size)

Spherical: $\Sigma_i = \sigma_i^2 I$ (\hookrightarrow close is aligned)

Diagonal: $\Sigma_i = \text{diag}(\sigma_{i,1}^2, \dots, \sigma_{i,d}^2)$ (\hookrightarrow aligned)

Tied: $\Sigma_i = \dots = \Sigma_k$ (\hookrightarrow same shape and size)

Full: Σ arbitrary (\hookrightarrow linear dbl)

K-means:

Special case of Hard-EM

• uniform weights $W_{i,k} = \frac{1}{k}$ and identical,

Spherical covariances $\Sigma_1 = \dots = \Sigma_k = \sigma^2 I$

Limiting case of Soft-EM

+ variances tending to 0: $\sigma \rightarrow 0$

Initialization

Weights: uniform distribution

Means: randomly, K-means++

Variances: spherical, e.g. according to empirical variance in data

K: cross-validation (only works if data truly gen. from a GMF)

Degeneracy

Problem: loss function not lower bounded:

\hookrightarrow loss converge to -infinity as $\mu \rightarrow \infty$ and $\sigma \rightarrow 0$

\hookrightarrow "Optimal" GMM chooses $k=n \rightarrow$ Overfitting

Solution:

- lower bound the width of the Gaussians

\hookrightarrow add $\gamma^2 I$ to variance Σ_j (regularizer)

Gaussian-Mixture Bayes Classifier

Estimate class prior $p(y)$

Estimate cond. distr. for each class as GMM:

$$P(x|y) = \sum_{j=1}^k W_j^{(y)} N(x; \mu_j^{(y)}, \Sigma_j^{(y)})$$

GAUSSIAN MIXTURE M₀

Assume data is generated from a convex-combination of Gaussian distributions:

$$p(x|\theta) = P(X|\mu, \Sigma, \omega) = \sum_{j=1}^k \omega_j N(x; \mu_j, \Sigma_j)$$

Objective (non-convex, not closed form) \hookrightarrow x_i are iid

$$L(\mu, \Sigma, \omega) = -\sum_{i=1}^n \log \sum_{j=1}^k \omega_j N(x_i; \mu_j, \Sigma_j)$$

Fitting a GMM = training a GBC without labels

Hard-EM

E-step: Predict most likely class for each data pts.

$$\zeta_i = \underset{j}{\operatorname{argmax}} P(z_i | x_i, \theta^{(t-1)})$$

$$= \underset{j}{\operatorname{argmax}} P(z_i | x_i, \theta^{(t-1)}) \cdot P(x_i | z_i, \theta^{(t-1)})$$

M-step: Compute MLE as for GBC (closed form)

$$\theta^{(t)} = \underset{\theta}{\operatorname{argmax}} P(D|\theta)$$

Problems: assigns fixed labels, even when model is uncertain

\hookrightarrow works poorly if clusters are overlapping

Expectation maximization

E-step: Calculate cluster membership weights for each pt.

$$y_j(x) = P(z=j | x, \Sigma, \mu, \omega) = \frac{P(x=x | z=j) \cdot P(z=j)}{P(x=x)}$$

$$= \frac{w_j \cdot P(x | \Sigma_j, \mu_j)}{\sum_{i=1}^k w_i \cdot P(x | \Sigma_i, \mu_i)}$$

$$y_j(x_i) = [\underbrace{j}_{=y_i}] \quad (\text{SSL})$$

M-step: Fit clusters to weighted data pts (closed form MLE)

$$w_j^{(t)} = \frac{1}{n} \sum_{i=1}^n y_j^{(t)}(x_i), \quad \mu_j^{(t)} = \frac{\sum_{i=1}^n y_j^{(t)}(x_i) \cdot x_i}{\sum_{i=1}^n y_j^{(t)}(x_i)}$$

$$\Sigma_j^{(t)} = \frac{\sum_{i=1}^n y_j^{(t)}(x_i) \cdot (x_i - \mu_j^{(t)}) \cdot (x_i - \mu_j^{(t)})^T}{\sum_{i=1}^n y_j^{(t)}(x_i)}$$

LLMs

Joint distr. over dependent categorical rand. var.

$$P(x_t=x | x_{1:t-1}=x_{1:t-1}) \approx P(x_t=x | x_{t-k:t-1}=x_{t-k:t-1}, \theta)$$

$$:= \text{Cat}(x_t) \text{ softmax}(\beta_t(x_{t-k:t-1}, \theta))$$

Transformer model

Input: $Z_o = XW_e + W_p$

Output: W_p : fixed pos. emb.

Trans. block: $Z_i = \text{transformer_block}(Z_{i-1}) \forall i \in \{1, \dots, n\}$

Output: $P = \text{softmax}(Z_n W_e^T)$

Transformer block: "masked multi-head self attention"

Attention mechanism learns to predict a weighted directed graph

- let words borrow features of words it attends to

Self-attention: restricts the nodes on both sides to be the same

Masking: only allows to attend to previous words

$$Z_i^{t+1} = \sum_{j=1}^n \text{score}_{i,j} Z_j^t; \quad \text{score}_{i,j} \propto \exp\left(\frac{g_{i,j}^T}{\sqrt{k}} + m_j\right)$$

$$m_{i,j} = \begin{cases} -\infty & \text{if } j > i \\ 0 & \text{otherwise} \end{cases}$$

Norm:

$$\|W\|_2 \leq \|W\|_1 \leq \sqrt{\|W\|_0 \cdot \|W\|_2}$$

Matrix:

$$(AB)^{-1} = B^{-1} A^{-1}, \quad (AB)^T = B^T A^T$$

$$(A^T)^{-1} = (A^{-1})^T, \quad \text{Tr}(AB) = \text{Tr}(BA)$$

Linear approximation of loss around point w^{new}

$$L(w^{\text{new}}) = L(w^{\text{new}} + \eta v) \approx L(w^{\text{new}}) + \eta \langle \nabla L(w^{\text{new}}), v \rangle$$

Convexity

$\alpha f + \beta g$; $\alpha, \beta \geq 0$, convex if f, g convex

$f \circ g$, convex if f convex and g affine or f non-dec and g convex

$\max(f, g)$, convex if f, g convex

GAUSSIAN MIXTURE M₀

Imputing missing features

Covariance matrix

$$\Sigma = \begin{bmatrix} \text{Var}[x_1] & \text{Cov}[x_1, x_2] \\ \text{Cov}[x_1, x_2] & \text{Var}[x_2] \end{bmatrix}$$

$$\mathbb{E}[x_1 | x_2 = x_2] = \mu_1 + \frac{\text{Cov}[x_1, x_2]}{\text{Var}[x_2]} (x_2 - \mu_2)$$

Convergence of EM

Monotonically increases the likelihood \rightarrow loc. optim.

\hookrightarrow initialization matter \rightarrow rerun and use sol. with largest LR

General EM

E-step: Calculate expected complete data log-likelihood

$$Q(\theta; \theta^{(t-1)}) = \mathbb{E}_{z_{1:n}} [\log P(x_{1:n}, z_{1:n} | \theta^{(t-1)}) | x_{1:n}, \theta^{(t-1)}]$$

$$= \sum_{i=1}^n \sum_{j=1}^k P(z_i=j | x_i, \theta^{(t-1)}) \cdot \log P(x_i, z_i | \theta^{(t-1)})$$

M-step: Maximize

$$\theta^{(t)} = \underset{\theta}{\operatorname{argmax}} Q(\theta; \theta^{(t-1)})$$

KL divergence

$$KL(p||q) = \sum_x p(x) \cdot \frac{p(x)}{q(x)} \geq 0, \quad KL(p||q)=0 \text{ if } p=q$$

Distributions

$$\text{Gaussian: } f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$\text{Poisson: } f(k; \lambda) = \frac{\lambda^k e^{-\lambda}}{k!}; \quad P(X \leq \lambda) = \sum_{k=0}^{\lambda} \frac{\lambda^k e^{-\lambda}}{k!}$$

Cont. prob: $P(Z_i) = P(z_i | Z_i \geq \lambda) \cdot P(Z_i \geq \lambda)$

Expectation and variance

$$\text{Var}[X] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})^2$$

$$\text{Var}[ax+b] = a^2 \text{Var}[X]$$

$$\text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y] + 2 \cdot \text{cov}(X, Y)$$

$$\mathbb{E}[X \cdot Y] = \mathbb{E}[X] \cdot \mathbb{E}[Y]$$

$$\text{Cov}(X, Y) = \frac{1}{N-1} \sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})$$

Bayes Theorem

$$P(y|x) = \frac{1}{P(x)} \cdot \frac{p(y) \cdot p(x|y)}{\sum_y p(y) \cdot p(x|y)}$$

Linear alg:

X has full rank $\Rightarrow X^T X$ invertible

$\Rightarrow X^T X$ is pd \Rightarrow $L(w) = \frac{1}{2} \|y - Xw\|^2$ is strongly convex

XX^T is always at least psd!

Inverse

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}, \quad A^{-1} = \det(A) \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

A diagonal matrix is invertible only if all its diagonal elements are non-zero

Positive (semi-) definite matrices

M is psd $\Rightarrow \forall x \in \mathbb{R}^n: x^T M x \geq 0 \Rightarrow \lambda_i \geq 0$

pd if and only if its trace AND determinant are positive > 0

Eigenvectors and eigenvalues

$$AV = \lambda V \Leftrightarrow A^2 V = \frac{1}{\lambda} V$$

SVD:

$$X$$