

Supervised learning learn model using collection of training e.g.s.

Regression predict a scalar value

Classification predict a discrete class target

K-NEAREST NEIGHBOURS non-parametric

classify a new input x based on its k nearest neighbours

Decision Boundary: boundary between regions of input space associated w/ diff. categories

Hyperparameter k controls # of neighbours

↓ good at capturing fine-grains, may overfit

↑ stable predictions, may underfit

Curse of Dimensionality ↑ dimensionality, # of data points needed for good performance ↑ exponentially

Sensitive to data scales → just normalize stuff

More memory/computation cost

DECISION TREES parametric

Make predictions by splitting features according to a tree

Decision boundary: axis aligned planes

Greedy algorithm based on picking features that maximize information gain

ENTROPY

$$H(Y) = -\sum_{y \in Y} p(y) \log_2 p(y)$$

$$H(X, Y) = -\sum_{x \in X} \sum_{y \in Y} p(x, y) H(Y | X=x)$$

$$H(Y|X) = \sum_{x \in X} p(x) H(Y | X=x) = -\sum_x \sum_y p(x, y) \log_2 p(y | x)$$

$$IG(Y|X) = H(Y) - H(Y|X)$$

LINEAR REGRESSION

Linear functions of features $y = w^T x + b$

b often added into w by padding x with 1's

Loss function $\mathcal{L}(y, t)$ used for learning

Cost function $J(w, b) = \mathcal{L}(y, t)$ averaged

Can vectorize inputs to abuse GPU

Minimize avg. training error via direct soln, GD methods

Feature mapping to put features on a diff. linearly separable space

GRADIENT DESCENT

$$w_j \leftarrow w_j - \alpha \frac{\partial}{\partial w_j} \quad w \leftarrow w - \alpha \sum_{i=1}^N (y^{(i)} - t^{(i)}) x_i^{(i)}$$

Stochastic update weights using just one training example

Mini batch split dataset into mini-batches, average each batch

If **Convex**, all critical points are global optimum

$$\hookrightarrow f((1-\lambda)x_0 + \lambda x_1) \leq (1-\lambda)f(x_0) + \lambda f(x_1)$$

$$RL \quad Q^k(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^k R_{t+k} \mid s, a \right] = r(s, a) + \gamma \sum_s P(s' \mid s, a) \sum_a \pi(a \mid s) Q^k(s', a')$$

$$Q^*(s, a) = r(s, a) + \gamma \sum_s P(s' \mid s, a) \max_{a'} Q^k(s', a')$$

$$V^k(s) = \mathbb{E}_{\pi} [G_k \mid s] = \sum_a \pi(a \mid s) [r(s, a) + \gamma \sum_s P(s' \mid s, a) V^k(s')]$$

REGULARIZATION

Adding penalties (e.g. L^2 penalty $R(w) = \frac{1}{2} \|w\|_2^2$)

$$\hookrightarrow J_{reg}(w) = J(w) + \lambda R(w), \text{ so } w \leftarrow (1-\alpha\lambda)w - \alpha \frac{\partial J}{\partial w}$$

LINEAR CLASSIFICATION

Linear functions of features $z = w^T x + b$ along w/ decision rule $y = \mathbb{I}[z \geq 0][1] + \mathbb{I}[z < 0][0]$

if data is linearly separable, we can find decision boundary

Logistic regression uses $y = \sigma(z)$

Loss function 0-1: ∇ issues, Squared loss: blows up extremes

[Logistic + SE]: makes extreme misclassifications look normal

$$\text{Binary Cross Entropy } \mathcal{L}_{CE} = -t \log y - (1-t) \log (1-y)$$

Closed-form solution

MULTI-CLASS CLASSIFICATION

Targets belong to discrete set, one hot encoding \mathbb{E}^K
 e^{z_k}
 $z = w^T x + b, \quad y_k = \text{Softmax}(z_1, \dots, z_K) = \frac{e^{z_k}}{\sum_k e^{z_k}}$

$$\text{Loss function } \mathcal{L}(y, t) = -\sum_k t_k \log y_k = -t^T (\log y)$$

NEURAL NETWORKS parametric

$$y = f^L \circ \dots \circ f^1(x), \quad f^{(i)} = \phi(wx + b)$$

Hidden unit, layer, weight, activation function

Each layer is a feature detector

Expressivity universal func. approximators (non-lin. act. func.)

Regularize by early stopping

Backpropagation learning algorithm ($\bar{y} = \frac{\partial \mathcal{L}}{\partial y}$ error signal)

Convolutional NN each set of hidden units looks at local region

weights shared between all image locations equivariant

Convolution detection, filters, produces feature maps

Pooling reduces size of rep., builds invariance, max selects max translations not applied

MODEL COMPLEXITY

Bias inability for method to capture true relationship, underfitting

Boosting merge different types of predictions, ↓ bias

Variance difference in fits between datasets, overfitting

$$\mathbb{E}[(y - t)^2 \mid x] = (y - \mathbb{E}[y])^2 + \text{bias} + \text{variance} + \text{Bayes error}$$

Bagging merge same type of predictions, ↓ variance

Bayes Optimality Bayes error due to inherent unpredictability of targets, achieving Bayes error is Bayes optimal

Decide H_i if $P(H_i \mid Z) > P(H_j \mid Z) \forall i \neq j$

Hyperparameter tuned using validation set, not learned

HYPOTHESIS SPACE

Hypothesis $f: X \rightarrow T$ from input to target space

If datasets were not naturally biased, all ML algos. would be the same

If $H_B \subseteq H_A$, model A is more expressive than B

Generally, the larger the more accurate.

PROBABILISTIC MODELS

Assuming outcomes $\{x_1, \dots, x_N\}$ are iid, log-likelihood of Bernoulli RV is $\ell(\theta) = \log \left(\prod_{i=1}^N \theta^{x_i} (1-\theta)^{1-x_i} \right)$

$$= \sum_{i=1}^N x_i \log \theta + (1-x_i) \log (1-\theta)$$

Maximum likelihood criteria pick $\hat{\theta} = \operatorname{argmax}_{\theta} \ell(\theta)$

Find optimal solutions by finding critical points, $\frac{d\ell}{d\theta} = 0$.

Discriminative classifiers feature $x \rightarrow$ class prob. $p(y|x)$

How do I learn classes from inputs?

Generative class $p(y) \rightarrow$ prob. of class given label $p(x|y)$

What does data for a class look like? Use Bayes' Rule

Naive Bayes model modes distribution of inputs for a label

Assumes features are cond. indep. given class

$$\begin{aligned} p(c|x) &= \frac{p(x|c) p(c)}{\sum_c p(c) p(x|c)} && \text{Prior} \\ \text{posterior for class} \\ l(\theta) &= \sum_{i=1}^N \log p(c^{(i)}, x^{(i)}) = \sum_{i=1}^N \log(p(x^{(i)}|c^{(i)}) p(c^{(i)})) \\ &= \sum_{i=1}^N \log p(c^{(i)}) + \sum_{j=1}^D \sum_{i=1}^N \log p(x_j^{(i)} | c^{(i)}) \\ &\quad \text{labels, prior} \quad \text{feature } x_j, \text{ posterior} \end{aligned}$$

Beta distribution $p(\theta; a, b) = \frac{(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}, E[\theta] = \frac{a}{a+b}$

Posterior distribution $p(D|\theta) = \frac{p(\theta) p(D|\theta)}{\int p(\theta) p(D|\theta) d\theta}$

Maximum A-Posteriori $\hat{\theta}_{MAP} = \operatorname{argmax}_{\theta} \log p(\theta) + \log p(D|\theta)$

MULTIVARIATE GAUSSIAN (normal) \uparrow posterior probability

(μ, Σ) uniquely define a multivariate Gaussian distribution

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

$$\hat{x} = Sx + b, \quad E[\hat{x}] = b, \quad \text{Cov}[\hat{x}] = SS^T$$

Since $\Sigma = E[(x - \mu)^T (x - \mu)]$ is symmetric, $\Sigma = Q \Lambda Q^T$

$$\begin{aligned} l(\mu, \Sigma) &= \sum_{i=1}^N -\log(2\pi)^{d/2} - \log |\Sigma|^{1/2} - \frac{1}{2} (x^{(i)} - \mu)^T \Sigma^{-1} (x^{(i)} - \mu) \\ \hat{\mu} &= \frac{1}{N} \sum_{i=1}^N x^{(i)} \quad \hat{\Sigma} = \frac{1}{N} (X - 1\mu^T)^T (X - 1\mu^T)^T \end{aligned}$$

Linear Regression $t|x \sim N(w^T x, \sigma^2)$, max log-likelihood

$$\frac{1}{N} \sum_{i=1}^N \log p(t^{(i)} | x^{(i)}; w, b) = \dots = \text{const.} \cdot \frac{1}{2N\sigma^2} \sum_{i=1}^N (t^{(i)} - w^T x^{(i)})^2$$

L_2 regularization MAP w/ Gaussian prior, $w \sim N(m, S)$

$$\log p(w) = \frac{1}{2} (w - m)^T S^{-1} (w - m) + \text{const.} = -\frac{1}{2\sigma^2} \|w\|^2 + \text{const.}$$

GAUSSIAN DISCRIMINANT ANALYSIS

Assumes $p(x|t)$ distributed according to multivar Gaussian.

$$\text{assuming binary class: } \phi = \frac{1}{N} \sum_{i=1}^N r_i^{(i)} \quad \mu_k = \frac{\sum_{i=1}^N r_k^{(i)} \cdot x^{(i)}}{\sum_{i=1}^N r_k^{(i)}} \quad \Sigma_k = \frac{1}{\sum_{i=1}^N r_k^{(i)}} \sum_{i=1}^N r_k^{(i)} \|x^{(i)} - \mu_k\|^2$$

Decision boundary

$$\log p(t_k|x) = -\frac{D}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log p(t_k) - \log p(x)$$

$$\Sigma_k = \Sigma_1 \text{ then linear boundary } (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) = (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \text{const.} \quad \text{conic}$$

Gaussian Naive Bayes

$$\mu_{jk} = \frac{\sum r_k^{(i)} x_j^{(i)}}{\sum r_k^{(i)}} \quad \sigma_{jk}^2 = \frac{\sum r_k^{(i)} (x_j^{(i)} - \mu_{jk})^2}{\sum r_k^{(i)}} \quad r_k^{(i)} = \mathbb{I}[t^{(i)} = k]$$

Isotropic?

If $\Sigma = \sigma^2 I$, we only need one parameter

$$\log p(t_k|x) - \log p(t_l|x) = -\frac{1}{2\sigma^2} [\|x - \mu_k\|^2 - \|x - \mu_l\|^2]$$

PRINCIPAL COMPONENT ANALYSIS

Dimensionality Reduction map to lower dim. space to retain only useful info, ideally close to intrinsic dimension

PCA is unsupervised, linear dimensionality reduction, projections

$$\text{Proj}_S(x) = \sum_{i=1}^k z_i u_i \text{ where } z_i = x^T u_i, \quad \text{Proj}_S(x) = Uz \text{ where } z = U^T x$$

$\hat{x} = Uz + \hat{\mu}$ can just shift to arbitrary origin.

$$\text{Minimize reconstruction error } \min_u \frac{1}{N} \sum \|x^{(i)} - \hat{x}^{(i)}\|^2$$

$$\text{Maximize variance of reconstructions } \max_u \frac{1}{N} \sum \|\hat{x}^{(i)} - \hat{\mu}\|^2$$

$$\text{Empirical Cov}(x) = \hat{\Sigma} = \sum_{i=1}^N (x^{(i)} - \hat{\mu})(x^{(i)} - \hat{\mu})^T \text{ symmetric, PSD}$$

Principal components eigenvectors, optimal subspace spanned by top k .

$$\hookrightarrow k=1, \frac{1}{N} \sum \|\hat{x}^{(i)} - \hat{\mu}\|^2 = \dots = u^T \hat{\Sigma} u = u^T Q \Lambda Q^T u = \alpha^T \Lambda \alpha = \sum_{j=1}^D \lambda_j a_j^2$$

MATRIX FACTORIZATION

Assume $x^{(i)} \approx \tilde{x}^{(i)} = Uz^{(i)}$ centered at $\hat{\mu} = 0$.

$$\text{Squared reconstruction error } \sum_{i=1}^N \|x^{(i)} - Uz^{(i)}\|^2 = \|X - UU^T\|_F^2$$

Sparse data matrix $R \approx UU^T$, need to fill in values w/ alternating least squares

$$\text{Squared error loss } \min_{U, Z} \frac{1}{2} \sum_{i,j} (R_{ij} - U_{ij} z_j)^2 = \sum_{j: R_{ij} \neq 0} (R_{ij} - U_{ij} z_j)^2$$

$$\text{Minimize wrt. } u; \text{ first, } u_i = \left(\sum_{j: R_{ij} \neq 0} z_j \right)^{-1} \sum_{j: R_{ij} \neq 0} R_{ij} z_j, \text{ then } z_j$$

AUTOENCODER

Feed-forward NN to take input x & predict x

Encoder (hopefully removes noise), **Decoder** recreates data

1 hidden layer linear reduce dimensions for code vector, 2 loss

Nonlinear learns to map to a manifold, more powerful

PCA, matrix factor., auto. are latent variable models

Assume data depends on some never used latent variables.

K-MEANS

(multimodal)

Clustering grouping data points in clusters w/ no labels

Find cluster centers m & assignments r to min. sum of 2 distances

$$\min_{m_1, m_2, r_1, r_2} J(\{m_1, m_2\}, \{r_1, r_2\}) = \min_{m_1, m_2} \sum_{n=1}^N \sum_{k=1}^2 r_k^{(n)} \|m_k - x^{(n)}\|^2$$

$$\text{Assignment } r = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_j \|m_j - x^{(n)}\|^2 \\ 0 & \text{otherwise} \end{cases} \quad \text{Refitting } m_k = \frac{\sum r_k^{(n)} x^{(n)}}{\sum r_k^{(n)}}$$

Alternating minimization, block coordinate descent

$$\text{Soft } r_k^{(n)} = \text{softmax}(-\beta \{ \|m_k - x^{(n)}\|^2 \}_{k=1}^K), \quad \beta \rightarrow \infty \text{ then } \rightarrow \text{normal k-means}$$

GAUSSIAN MIXTURE MODELS

$p(z)$ in $p(x) = \sum_z p(x|z) = \sum_z p(x|z) p(z)$ categorical \rightarrow mixture model

GMM $p(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$, fit params. w/ max likelihood.

EXPECTATION-MAXIMIZATION

Find (local) max likelihood params. of model

$$\log p(x|\theta) = \sum_{i=1}^N \log \sum_{k=1}^K p(x^{(i)} | z^{(i)}; \{\mu_k\}, \{\Sigma_k\}) p(z^{(i)} | \Pi)$$

E-step compute post. prob. over z , given θ compute $z^{(i)}$

$$\text{Assign responsibility } r_k^{(i)} = \Pr(z^{(i)} = k | x^{(i)}; \theta)$$

M-step given $z^{(i)}$ learn θ , apply max likelihood updates

Gaussian discriminant analysis