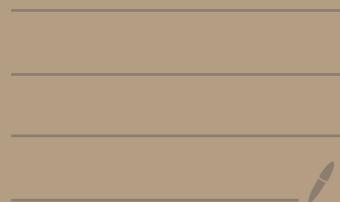


CS335/337 : AI & ML

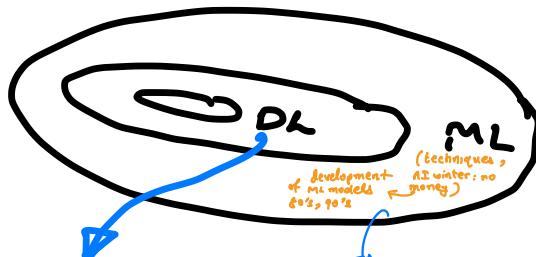
(Prof. Preeti Jyothi)



learn: Make accurate predictions or decisions based on data by optimizing model
models are lossy by nature.

"All models are wrong but some are useful"

Performance measure to identify which f is better



Origin:
Perceptrons
(lin. models)
AI (1950s: perceptrons)

NNs,
deep learning
models

SVMs, MCMC,
hidden markov, CNNs

ML & Statistics

proving almp-
-otic bounds

ML \subseteq CS \rightarrow scalability,
heuristics,
efficient
algorithms,
finite sample

When ML?

1. For tasks that are easy for humans but are complex for computer systems to emulate (hard to code)

- Computer vision (solved problem)
- Natural language (translation, Q/A, identifying sentiment)
- Speech (recognize speech, speaking sentences naturally)
- Game playing (play games like chess, Go, DOTA, Poker, etc.) {RL?}

- Robotics (Walking, jumping, displaying emotions) { usually not mh }
 - Self driving, navigating maze, etc.
2. For tasks that are beyond human capabilities
- Analysis of large and complex datasets
 - IBM Watson's Jeopardy machine (parse nat. lang.), health analyst (symptoms, history \rightarrow medical diagnosis)

Kaggle : Community to solve various problems .

↳ Survey '17 - '21 : { Logistic reg > decision trees > models to be covered in grad boosting > CNNs > Bayesian > this course . Dense NNs > RNNs }

1) Supervised learning : Decision trees, neural networks, etc.

Let D be a dataset with n training examples,

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

$$\text{(or) } D = \{(x_i, y_i)\}_{i=1}^n, \text{ where } x_i \in \mathcal{X},$$

\mathcal{X} is feature space. ($x_i \in \mathbb{R}^d$, x_i is a d -dimensional input real-valued vector)

$y_i \in \mathcal{Y}$, \mathcal{Y} is the label space

1. $\mathcal{Y} = \{0, 1\}$ or $\{-1, 1\}$, $|\mathcal{Y}| = 2$: Binary classification

2. $\mathcal{Y} = \{1, 2, \dots, k\}$, $k > 2$: Multiclass classification

3. $\mathcal{Y} = \mathbb{R}$: Regression

Objective : Let the instances in D be drawn from some unknown (true) probability distribution P over $\mathcal{X} \times \mathcal{Y}$.

Find an $h: \mathcal{X} \rightarrow \mathcal{Y}$ s.t. for a new test instance (x, y) , $h(x) = y$ with high probability
 $(h(x) \approx y)$

2) **Unsupervised Learning**: k-means, PCA, mixture models

$$\mathcal{D} = \{x_1, \dots, x_n\}$$

Identify structural properties to identify clusters, partition the data \mathcal{D}

3) **Reinforcement learning**: Not covered in course
(CS747, CS748)

$h \in \mathcal{H}$ (hypothesis class)

↓ * choice of model (e.g. linear, etc.)
what is the best model

Course logistics : refer moodle slides

Course project : 3-4 members

↳ pick doable problems, many datasets, many libraries

Lecture 2

(3/8/2023)

Preliminaries $\mathcal{X} \equiv \text{Input / Feature / Attribute space}$
 $y \equiv \text{Output / Label / Response}$

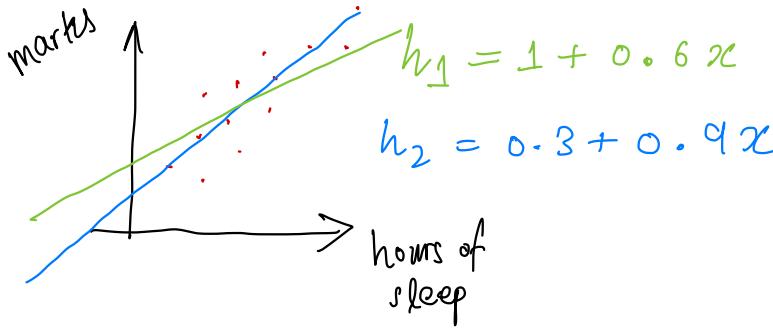
Linear Regression: $\mathcal{X} = \mathbb{R}^d$, $y = \mathbb{R}$

Consider a function $f: \mathcal{X} \rightarrow \mathcal{Y}$. Given a dataset \mathcal{D} which consists of n -tuples (x, y) (training dataset)
goal is to find a hypothesis h that ideally closely approximates f where, $h \in \mathcal{H}$ (\mathcal{H} is hypothesis class)

Three important questions

1. What predictors are permissible (characterization of $h \in \mathcal{H}$)? (Hypothesis class)
2. How do I measure if a predictor is good? (loss/error function)
3. How do we find the best predictor? (Optimization problem)

1. Possible predictors



$$h_w(x) = w_0 + w_1 x, \quad w = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$$

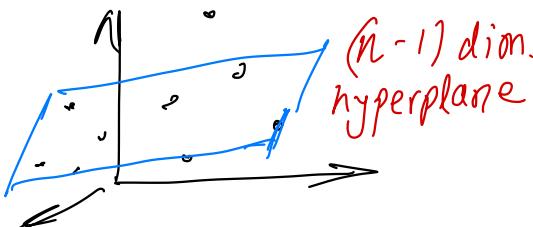
Vector notation $h_w(x) = w^T x \quad x = \begin{bmatrix} 1 \\ x \end{bmatrix}$

($h_w(x)$ is a linear combination of features)

even if $x = \begin{bmatrix} 1 \\ x^2 \end{bmatrix}$, model is still linear regression

$$h_w(x) = w^T x, \quad w = \begin{bmatrix} w_1 \\ \vdots \\ w_{d+1} \end{bmatrix}, \quad x = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$\mathcal{H} = \{h_w : w \in \mathbb{R}^{d+1}\}$ (hypothesis class)



2 • Quantify a measure of prediction

loss function
(error function)
(objective)

$L(h, D)$
↑ hypothesis ↑ dataset

In case of linear regression, we can work with $L(w, D)$

Linear Regression

$$L(w, D_{\text{test}}) = (y - h_w(x))^2$$

$|y - \hat{y}|$ referred to as residuals

least squares loss

$D_{\text{test}} = \{(x, y)\}$
(test instance)

for absolute loss,

1. For outliers, it's more forgiving since grows slower than square loss
2. Analytic optimum issues due to non-diff.

$$L(w, D_{\text{train}}) = \frac{1}{n} \sum_{i=1}^n (y_i - w^\top x_i)^2$$

omitted at times
for simplicity

Now, we stack x_i 's to make

$D_{\text{train}} = \{(x_i, y_i)\}_{i=1}^n$

$$X = \begin{bmatrix} x_1^\top \\ \vdots \\ x_n^\top \end{bmatrix}_{n \times (d+1)}, \quad Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad w = \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix}_{(d+1) \times 1}$$

$$L(w, D_{\text{train}}) = \frac{1}{n} \| \gamma - Xw \|_2^2 \quad (\text{l2 norm squared})$$

3. Finding the best predictor

$$h^* = \arg \min_h L(h, D_{\text{train}})$$

(training loss)

$$w^* = \arg \min_w L(w, D_{\text{train}})$$

$$w_{\text{LS}}^* = \arg \min_w \sum_{i=1}^n (y_i - w^\top x_i)^2$$

Least squares solution

Closed form solution for 3-dimensional data

$$w^* = \arg \min_w \underbrace{\sum_i (y_i - w^\top x_i)^2}_{L_w}$$

Take the derivative of L_w w.r.t w , set to 0
and solve for w

$$\nabla_w L = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial w_n} \end{bmatrix} = 2 \left(\sum_i (y_i - w^\top x_i) x_i \right) = 0$$

$$\Rightarrow \sum_i y_i x_i - \sum_i (w^\top x_i) x_i = 0$$

$$\Rightarrow w = (X^\top X)^{-1} X^\top \gamma$$

For data D , $h(x) = w_0 + w_1 x$

confirm that $w_1^* = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}$

$$\bar{x} = \frac{\sum_i x_i}{n},$$

$$\bar{y} = \frac{\sum_i y_i}{n}$$

Completing the linear regression proof,

$$\sum_i y_i x_i - \sum_i (w^T x_i) x_i = 0$$

$$\Rightarrow X^T y - X^T X w = 0 \quad X^T = [x_1 \dots x_n]$$

$$\Rightarrow w = (X^T X)^{-1} X^T y \quad y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

Note: If $X^T X$ is not inv, we get multiple pol's e.g. infinite planes pass through a line, happen when cols of X are h.D. i.e. are "collinear" or $(n-2)$ hyperplane lying.

Lecture 3

(8/8/28)

1. Converting non-linear regression to linear regression

In the general regression case, h need not be linear. Consider the following simple case: $h_w = w_0 + w_1 x + w_2 x^2$

$\phi: \mathbb{R} \rightarrow \mathbb{R}^3$ by $\phi(x) = \begin{pmatrix} 1 \\ x \\ x^2 \end{pmatrix}$, $h_w(x) = w^\top \phi(x)$

$$w = (\phi^\top \phi)^{-1} \phi^\top y, \text{ where } \phi = \begin{pmatrix} \phi(x_1)^\top \\ \phi(x_2)^\top \\ \vdots \\ \phi(x_n)^\top \end{pmatrix}$$

(General case)

$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$. $\phi(x) = \begin{pmatrix} \phi_1(x) \\ \dots \\ \phi_D(x) \end{pmatrix}$ ($x_i \in \mathbb{R}^d$)

$$h_w(x) = w^\top \phi(x).$$

$$h(w, D) = \frac{1}{n} \sum_{i=1}^n (h_w(x_i) - y_i)^2 = \frac{1}{n} \sum_{i=1}^n (w^\top \phi(x_i) - y_i)^2$$

So, we get normal solution,

$$w = (\phi^\top \phi)^{-1} \phi^\top Y$$

$$\text{where, } \phi = \begin{pmatrix} \phi(x_1)^\top \\ \dots \\ \phi(x_n)^\top \end{pmatrix}$$

Typical choices of basis transformations

1. Polynomial basis: $\phi(x) = (1, x, x^2, \dots, x^{D-1})^\top$ [the example we saw above]
2. Gaussian basis/Radial Basis functions: $\phi(x) = (1, e^{-\frac{(x-x_1)^2}{2\sigma_1^2}}, e^{-\frac{(x-x_2)^2}{2\sigma_2^2}}, \dots, e^{-\frac{(x-x_D)^2}{2\sigma_D^2}})^\top$
3. Fourier basis: $\phi(x) = (1, \cos(\omega_1^\top x), \sin(\omega_1^\top x), \dots, \cos(\omega_D^\top x), \sin(\omega_D^\top x))^\top$
4. Piecewise linear basis

5. Periodic basis

Some terminologies

- Training Data. The dataset $D = \{(x_i, y_i)\}_{i=1}^n$
 - Training Error. $h(w, D) = \sum_{i=1}^n (h_w(x_i) - y_i)^2$
 - Training - Dev - test $\approx 80:10:10$
 - Tuning hyperparameters: Grid Search, Random Search, Bayesian Opt.
- (training the model) (tuning the hyper-parameters) (eval- uate)

• Test error is a measure of the generalization of the model - better generalization is the holy grail of ML.

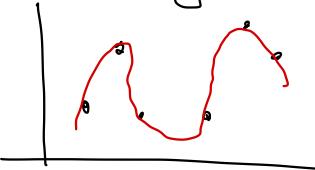
Underfitting and overfitting

$g_C = \text{all lines parallel to the axes e.g. } h_C(x) = C \text{ or } x_i = C$
 Too simple, cannot fit most data well - it is too simplistic and inflexible. Called "underfitting", where model is needed to be more complex.

Lecture 4

(08/08/23)

Overfitting



Training error: 0

Test error: high

How do we combat this?

1. **Coarse strategy**: tune k (order of polynomial) and evaluate on a dev set to pick a "good value".
2. **Regularization**: modify the loss function to explicitly constrain model complexity

$h_{\text{reg}}(\omega, D_{\text{train}})$

$$= L_{\text{MSE}}(\omega, D_{\text{train}}) + \lambda R(\omega), \text{ where } \lambda \geq 0.$$

$R(\omega)$: Regularizer that penalizes or shrinks the weights in ω

1st term: $h_{\text{reg}}(\omega, D_{\text{train}})$ = measure of fit to training data

2nd term: $R(\omega)$ = measure of model complexity

Penalty or shrinkage-based regularizations

1. L_2 reg.: $R(\omega) = \|\omega\|_2^2$ ("shrinks" norm)

] force ω to not become arbitrarily large.
 more variance in coeff measured by $\|\omega\|$.

2. L_1 regularization: $R(\omega) = \|\omega\|_1$

Induced regressions

1. Ridge regression (L_2 -normalised regression)

$$\omega_{\text{RIDGE}} = \arg \min_{\omega} \|y - \phi\omega\|^2 + \lambda \|\omega\|_2^2$$

Equivalent to

$$\omega_{\text{RIDGE}} = \arg \min_{\omega} \|y - \phi\omega\|^2 \quad \text{s.t. } \|\omega\|_2^2 \leq t$$

Note: It is common to leave out ω_0 (intercept) in ridge regression since, indep. of x_i , we might miss out on a best fit

Note: we dropped the y_n for $\|y - \phi\omega\|^2$ hence, λ must be scaled appropriately

Solving,

$$\nabla L_{\text{RIDGE}}(\omega) = 0 \Rightarrow -2\phi^T(y - \phi\omega) + 2\lambda\omega = 0$$

$$\omega_{\text{RIDGE}} = (\phi^T\phi + \lambda I)^{-1}\phi^T y$$

may not be invertible \uparrow mult be invertible! $\left\{ \begin{array}{l} \text{since } \phi^T\phi \text{ is p.c.d.} \\ \text{so, } \phi^T\phi + \lambda I, \lambda > 0, \text{ is p-d} \end{array} \right\}$

2. Lasso Regression (L_1 -regularised regression)

$$w_{\text{Lasso}} = \arg \min_w \|y - \phi w\|_2^2 + \lambda \|w\|_1$$

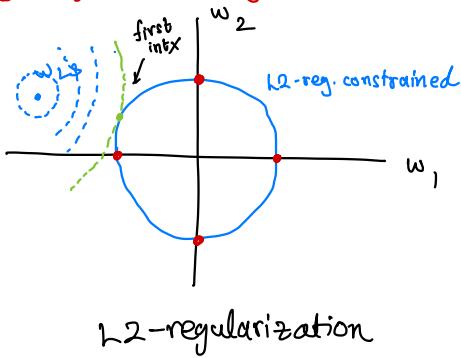
(no closed form, but other ways to solve)

Solve for w_{Lasso} using

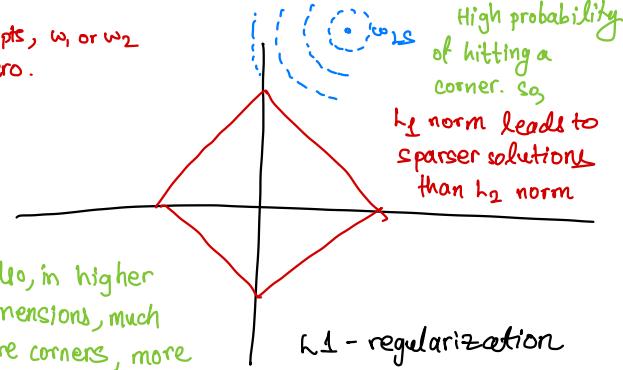
1. Quadratic Programming

2. Iterative optimization algorithm like "gradient descent"

L_1 regularization yields sparse weight vectors compared to L_2 regularization



At red pts, w_1 or w_2 is zero.



Also, in higher dimensions, much more corners, more chance of hitting them

L_1 -regularization

- When your data has outliers, sparsity is useful.

lecture 5

(10/08/2023)

Gradient descent (GD)

GD is an iterative optimization algorithm

General template:

Gradient descent

1. Initialize w (e.g. $w=0$)

$w \leftarrow w_0$ (e.g. 0, etc.)

[for iteration $t=0, 1, \dots$]

Direction: $-\nabla L(w_t)$

Repeat 1. choose a descent direction

$\alpha > 0$ (hyperparameter) (learning rate)

2. choose a step size

$w_{t+1} \leftarrow w_t - \alpha \nabla L(w_t)$

3. Update w

Common stopping criterion

Until stopping criterion is met

$\|\nabla L(w_t)\|_2 \leq \epsilon, \|w_{t+1} - w_t\| \leq \epsilon$

Gradient descent comes from

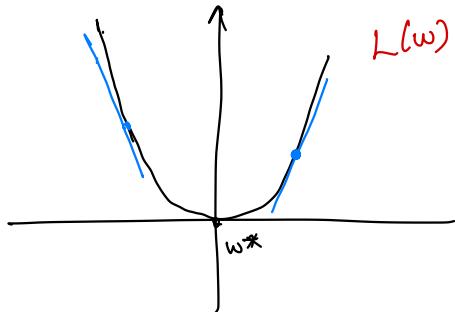
ϵ is hyperparameter

$$\text{Gradient } \nabla L(w) = \begin{bmatrix} \frac{\partial L(w)}{\partial w_1} \\ \vdots \\ \frac{\partial L(w)}{\partial w_n} \end{bmatrix}$$

gives you the direction of fastest increase.

descent update w in the direction of fastest decrease in L

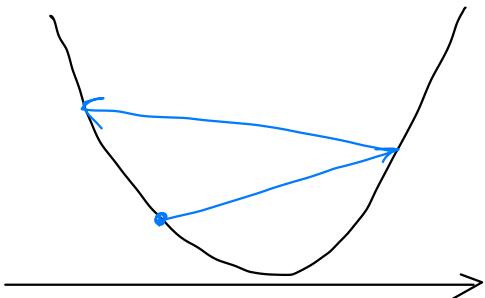
For 1D case and linear regression



what if step size is too small?
Takes longer to converge.

Step size determines how fast you converge.

what if step size is very large?



Loss can diverge in worst case,
GD will not converge

What is the weight update rule
in GD for linear regression?

$$w \leftarrow w - \alpha \nabla L(w)$$

$$L(w) = \sum_i (y_i - w^T x_i)^2$$

$$\begin{aligned} \nabla L(w) &= -2 \sum_i (y_i - w^T x_i) x_i \\ &= -2 \sum_i (y_i - \hat{y}_i) x_i \end{aligned}$$

Substituting, we get

$$w \leftarrow w - \alpha \sum_i (\hat{y}_i - y_i) x_i$$

To counter SGD issues,

Mini-batch gradient descent :

$$w \leftarrow w - \alpha \nabla L(w, D_{\text{batch}})$$

where $D_{\text{batch}} = \{(x_i, y_i)\}_{i=1}^B$, B = batch size

Batch size dictated by RAM/GPU, basically,
how many examples you can compute upon.

$$\text{GD : } w \leftarrow w - \alpha \nabla L(w, D_{\text{train}})$$

SGD (Stochastic Gradient descent)

$$w \leftarrow w - \alpha \nabla L(w, D_{\text{random}})$$

where, $D_{\text{random}} = \{(x_i, y_i)\}$,
 (x_i, y_i) is randomly sampled from D_{train}
• issue: if example picked is outlier,
not representative of D_{train} ,
we can take wrong steps, but
we have convergence guaranteed.

training becomes more stable
with larger batch sizes, typically
batch size = 32, 64.

(shuffle the dataset, then make
one pass over all examples via
minibatches, one pass over all
examples is called an epoch)

Probabilistic view of Linear Regression

For training data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$. Let the target y_i have some noise defined as,

$$y_i = f(x_i) + \varepsilon_i, \varepsilon_i \sim N(0, \sigma^2)$$

ε_i 's are independent and identically distributed zero mean gaussians with same variance σ^2 . ($\text{Cov}(\varepsilon_i, \varepsilon_j) = 0 \forall i \neq j$)

For linear regression,

$$\Leftrightarrow y_i = w^T x_i + \varepsilon_i, \varepsilon_i \sim N(0, \sigma^2)$$

$$\Leftrightarrow y_i \sim N(w^T x_i, \sigma^2)$$

$$\Leftrightarrow P(y_i | x_i, w) = N(w^T x_i, \sigma^2)$$

\hookrightarrow Likelihood

For training data $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$,

$$P(y_1, y_2, \dots, y_n | x_1, \dots, x_n, w) = \prod_i P(y_i | x_i, w) \quad \begin{matrix} (y_i \text{ are conditionally} \\ \text{independent given } x_i) \end{matrix}$$

$$\log P(y_1, y_2, \dots, y_n | x_1, \dots, x_n, w) = \sum_i \log P(y_i | x_i, w) \quad (\text{log likelihood})$$

The probabilities can get really close to zero, so high chance of numerical underflow, log probabilities are much more well behaved.

We want to find a w which maximises this likelihood, which will reduce to mean squared error.

14/08/2023

Lecture 6

$$Y = w^T X + \varepsilon, \varepsilon \sim N(0, \sigma^2) \quad \begin{matrix} (\text{multivariate gaussian}) \end{matrix}$$

$$\Leftrightarrow Y \sim N(w^T X, \sigma^2) \quad (\text{examples stacked across columns})$$

For training dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$

$$P(y_1, y_2, \dots, y_n | x_1, \dots, x_n) = \prod_i P(y_i | x_i, w) \quad (\text{likelihood})$$

$$\log P(y_1, y_2, \dots, y_n | x_1, \dots, x_n) = \prod_i \log P(y_i | x_i, w) \quad (\text{log-likelihood})$$

Maximum likelihood estimation : Find parameters w that maximize the probability of the observed data.

$$w_{\text{MLE}} = \arg \max_w \sum_i \log P(y_i | x_i, w)$$

Motivating Example

You want to estimate the probability of a biased coin landing on heads (θ)
 say your data is N coin tosses, with N_H heads and N_T tails (i.i.d.).
 What is the MLE estimate of θ ?

$$P(D|\theta) = \theta^{N_H} (1-\theta)^{N_T}$$

$$\theta_{MLE} = \arg \max_{\theta} \log P(D|\theta) = \arg \max_{\theta} (N_H \log \theta + N_T \log(1-\theta))$$

$$\text{i.e. } \frac{N_H}{\theta} + \frac{N_T}{1-\theta} (-1) = 0 \Rightarrow (1-\theta)N_H = \theta \cdot N_T \\ \Rightarrow N_H = \theta \cdot (N_H + N_T) = \theta \cdot N \\ \Rightarrow \boxed{\theta = \frac{N_H}{N}}$$

\Leftrightarrow Say you have a coin with bias $p \in \{0.4, 0.6\}$

Data: 3 coin tosses, 2H, 1T

What is the MLE of p .

Ans $P(D|p) \propto p^2(1-p)$

at $p=0.4$, $(0.4)^2(0.6)$ Hence, $MLE = 0.6$

$p=0.6$, $(0.6)^2(0.4)$

MLE for linear regression

$$\omega_{MLE} = \arg \max_{\omega} \sum_{i=1}^n \log P(y_i|x_i, \omega)$$

$$= \arg \max_{\omega} \sum_{i=1}^n \log N(\omega^T x_i, \sigma^2)$$

$$= \arg \max_{\omega} \sum_{i=1}^n \log \left[\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y_i - \omega^T x_i)^2}{2\sigma^2}} \right]$$

$$= \arg \max_{\omega} \sum_{i=1}^n \left[C - \frac{(y_i - \omega^T x_i)^2}{2\sigma^2} \right]$$

$$\omega_{MLE} = \arg \min_{\omega} \sum_{i=1}^n (y_i - \omega^T x_i)^2 \quad (\text{least squares error solution})$$

Homework: If we have a Laplacian noise instead of Gaussian, we get L_1 norm
 i.e. least absolute error.

Note: The limitation of MLE is the tendency to overfit if you have
 limited data. (entirely data dependent).

Bayesian Parameter Estimation

In ML, observations are random variables, but parameters are not. In the
 Bayesian framework, parameters are also random variables.

\Rightarrow Underlying prior distribution which encodes beliefs about parameters before observing data.

Note: Prior is a pre-determined choice, parametrized by some vars, these parameters will be further hyperparameters.

Maximum A posteriori Estimate (MAP)

$$P(\theta|D) \propto P(D|\theta) P(\theta)$$

posterior likelihood prior

$$\theta_{\text{MAP}} = \arg \max_{\theta} P(\theta|D) = \arg \max_{\theta} P(D|\theta) P(\theta) = \arg \max_{\theta} \log P(D|\theta) + \log P(\theta)$$

(MLE)

Coin example (MAP estimate)

$$P(D|\theta) = \theta^{N_H} (1-\theta)^{N_T}$$

what is a good prior on θ ? Beta distribution is a good candidate for a prior

$$P(\theta, \alpha, \beta) = \frac{1}{C} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

We want to pick a prior so that posterior has a similar functional form as the prior

Conjugate Priors

For a likelihood $P(D|\theta)$ coming from a family of distributions d_1 , a prior is said to be conjugate prior (from a family d_2) if the posterior distribution also comes from family d_2 .

Coin example $B(N_H, N_T, \theta) = \theta^{N_H} (1-\theta)^{N_T}$

$$\text{Beta}(\theta, \alpha, \beta) = C \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

$$\begin{aligned} P(\theta|D) &\propto P(D|\theta) P(\theta) \\ &\propto \theta^{N_H+\alpha-1} (1-\theta)^{N_T+\beta-1} \\ &\propto \text{Beta}(\theta, N_H+\alpha, N_T+\beta) \end{aligned}$$

Extra reading: Conjugate priors drawn from exponential families usually e.g. for bernoulli, binomial likelihood, beta distribution works. for multinomial, dirichlet prior for gaussian, gaussian prior MLE ≡ unreg. linear regression MAP ≡ reg. linear regression

(17/08/2023)

Lecture 7

Recap: MAP estimation

$$\theta_{\text{MAP}} = \arg \max_{\theta} \underbrace{\log P(D|\theta) + \log P(\theta)}_{\text{Posterior probability distribution}}$$

Coin toss problem $B(n_H, n_T, \theta) = \theta^{n_H} (1-\theta)^{n_T}$

$$\text{Beta}(\theta, \alpha, \beta) = \frac{1}{C} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$
 [C is a normalisation constant]

$P(\theta|D) \propto \text{Beta}(\theta, n_H+\alpha, n_T+\beta)$ {beta distribution is a conjugate prior for the bernoulli/binomial likelihood}

$$\theta_{\text{MAP}} = \arg \max_{\theta} P(\theta|D) = \arg \max_{\theta} (n_H+\alpha-1) \log \theta + (n_T+\beta-1) \log (1-\theta)$$

$$\theta_{\text{MAP}} = \frac{n_H+\alpha-1}{n+\alpha+\beta-2}$$

$$\frac{n_H+\alpha-1}{n+\alpha+\beta-2}; \quad \leftarrow \text{(pseudocounts)}$$

Note: $\theta_{\text{MAP}} \rightarrow \theta_{\text{MLE}}$ as $n \rightarrow \infty$ (Bernstein-von Mises theorem)

MAP estimate for Linear Regression

Consider the following Gaussian prior over W ,

$$P(W) = N(0, \frac{1}{\lambda} I) = \left(\frac{\lambda}{2\pi}\right)^{d/2} \exp\left(-\frac{\lambda}{2} W^T W\right) \quad \text{follows}$$

$$\left[\underset{w \in \mathbb{R}^d}{N(\mu, \Sigma)} = \frac{1}{(2\pi)^{d/2}} |\Sigma|^{1/2} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right) \right]$$

$$W_{MAP} = \arg \max_w \log P(D|w) + \log P(w)$$

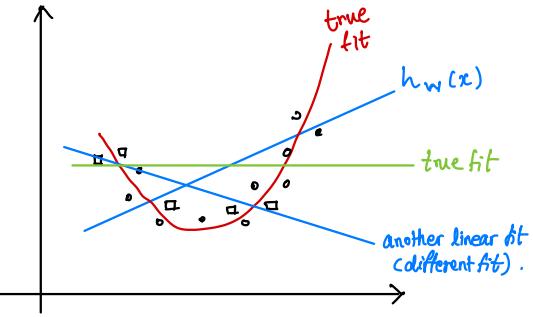
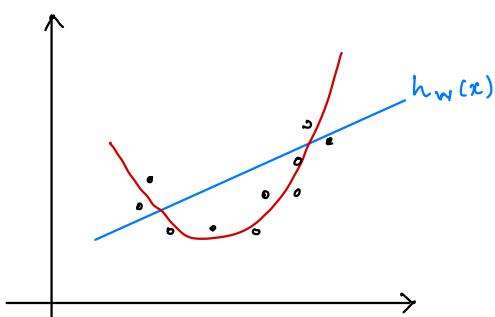
$$= \arg \max_w -\frac{1}{2\sigma^2} \sum_i (y_i - w^T x_i)^2 - \frac{\lambda}{2} \|w\|_2^2$$

$$= \arg \min_w \frac{1}{2\sigma^2} \sum_i (y_i - w^T x_i)^2 + \frac{\lambda}{2} \|w\|_2^2 \quad \{ \text{ridge or } L_2 \text{ regularised regression} \}$$

Note: If we use laplace prior instead of gaussian, we would get a L_1 regularisation.

$$(Laplacian) f(x| \mu, b) = \frac{1}{2b} \exp\left(-\frac{|x - \mu|}{b}\right)$$

Bias and Variance of Estimators



- How much does the mean of fit differ (linear, here) from true fit

← bias

inherent to choice of hypothesis class.

- What is the variance across fits when training set varies

← Variance

inherent to data

Bias - Variance Analysis for Linear Regression

Let $y = f(x) + \epsilon$, $\epsilon \sim N(0, \sigma^2)$. For a test point \tilde{x} , $\tilde{y} = f(\tilde{x}) + \tilde{\epsilon}$, $\tilde{\epsilon} \sim N(0, \sigma^2)$ the expected test error can be written as

$$\mathbb{E}_{\tilde{x}, \tilde{\epsilon}} [(\tilde{y} - h_w(\tilde{x}; \theta))^2]$$

shorthand for this, $\mathbb{E}[(\tilde{y} - h(\tilde{x}))^2]$

$$\mathbb{E}[(\tilde{y} - h(\tilde{x}))^2] = \mathbb{E}[\tilde{y}^2 - 2\tilde{y}h(\tilde{x}) + h(\tilde{x})^2] = \mathbb{E}[h(\tilde{x})^2] + \mathbb{E}[\tilde{y}^2] - 2\mathbb{E}[\tilde{y}]\mathbb{E}[h(\tilde{x})]$$

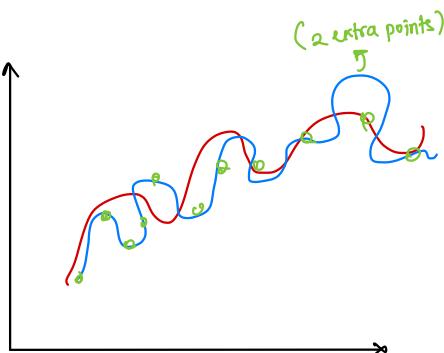
linearity of expectation

$$= \text{Var}(h(\tilde{x})) + \mathbb{E}[h(\tilde{x})]^2 + \dots$$

Note: \tilde{y} is true value, h is predictor so are independent R.V.s, h varies as it's a fn of dataset

$$\begin{aligned}
 &= \mathbb{E}[(h(\tilde{x}) - \bar{h}(\tilde{x}))^2] + \mathbb{E}[h(\tilde{x})]^2 + \mathbb{E}[(\tilde{y} - f(x))^2] + f(\tilde{x})^2 - 2\mathbb{E}[h(\tilde{x})]f(\tilde{x}) \\
 &= (\bar{h}(\tilde{x}) - f(\tilde{x}))^2 + \mathbb{E}[h(\tilde{x}) - \bar{h}(\tilde{x})^2] + \sigma^2 \\
 &= \boxed{\text{bias}^2 + \text{Variance} + \sigma^2}
 \end{aligned}$$

↳ putting $\tilde{y} = f(x) + \epsilon$



In high complexity models, variance is high, bias is low.

In low complexity (say lines || to x-axis), low variance, high bias.

Overfitting : High variance

Under-fitting : High bias

For ridge regression, if λ increases, variance decreases and bias increases (modelling capabilities reduced).

$$w_{\text{ridge}} = \arg \min_w \|Y - Xw\|^2 + \lambda \|w\|^2$$

Lecture

Logistic Regression

$f: X \rightarrow Y = \{0, 1\}$ (binary classification)

i.e. find the best deterministic estimate to a probability distribution over $x \in \{0, 1\}$. Loss f^n ?

$$1. \text{ Discrete loss (matrix } L\text{)}, \quad L_{ij} = \begin{cases} 1 & \text{if } j \neq i \\ 0 & \text{if } i=j \end{cases} \leftarrow \text{common}$$

$$2. \quad L_{ij} = |i-j| \quad (\text{absolute loss})$$

$$3. \quad L_{ij} = (i-j)^2 \quad (\text{squared loss})$$

$$\text{EPE}(f) = \mathbb{E}_x [\mathbb{E}_{Y|x} [L(Y, f(x))]]$$

minimizing pointwise, we have the regression

$$\begin{aligned}
 f(x) &= \arg \min_{y \in \mathcal{Y}} \mathbb{E}_{Y|x=x} [L(Y, y)] = \arg \min_{y \in \mathcal{Y}} \mathbb{E}_{Y|x=x} [\mathbb{1}_{[Y \neq y]}] \\
 &= \arg \max_{y \in \mathcal{Y}} P(Y=y | X=x)
 \end{aligned}$$

i.e. we classify based on most likely label (Bayes' classifier)

we have $y = \{0, 1\}$, so, $P(Y=1 | X=x) = 1 - P(Y=0 | X=x)$.

so, we have $f(x) = \begin{cases} 1 & [P(Y=1 | X=x) > 0.5] \\ 0 & [P(Y=1 | X=x) \leq 0.5] \end{cases}$

Logistic Regression

we attempt to approximate with a linear f^n .

$$f(x) = \begin{cases} 1, & \text{if } \beta^T x > 0.5 \\ 0, & \text{if } \beta^T x \leq 0.5 \end{cases}$$

but $\beta^T x$ can be negative. So, we use sigmoid f^n

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$f(x) = \begin{cases} 1, & \text{if } \sigma(\beta^T x) > 0.5 \\ 0, & \text{if } \sigma(\beta^T x) \leq 0.5 \end{cases}$$

ML estimate for β

we have loss for dataset $D = \{(x_i, y_i)\}_{i=1}^n$

$$L = \log \pi_i P(y_i | x_i, \beta)$$

$$\begin{aligned} \log y_i | x_i, \beta &= y_i \log P(y_i=1 | x_i, \beta) + (1-y_i) \log (1 - P(y_i=1 | x_i, \beta)) \\ &= y_i \log f(x) + (1-y_i) \log (1-f(x)) \end{aligned}$$

so, $\beta_{ML} = \arg \min -\sum_i y_i \log f(x_i) + (1-y_i) \log (1-f(x_i))$

The loss, $L = -\sum_i y_i \log f(x_i) + (1-y_i) \log (1-f(x_i))$ {binary cross-entropy loss}

$$L = -\sum_i y_i \log \sigma(\beta^T x_i) + (1-y_i) \log (1 - \sigma(\beta^T x_i))$$

$$\nabla_{\beta} L = \sum_i (\sigma(\beta^T x_i) - y_i) x_i$$

No closed form solution for β .

lecture

Logistic Regression

$$\begin{aligned} w^* &= \arg \max_w \sum_i \log P(y_i | x_i, w) \\ &= \arg \min_w \text{Cross-entropy-loss } (w, D) \end{aligned}$$

Regularized Logistic Regression

$$w_R^* = \arg \min_w L_{CE}(w, D) + \lambda \|w\|_2^2$$

In the unregularised case,

purple should get higher probability of 1 than green point (near decision boundary)

$$P(y|x, w) = \sigma(w^T x) = \frac{1}{1 + e^{-w^T x}}$$

Classifier will prefer high values of w to maximise $P(y|1, w)$.

But this will cause v. less difference b/w purple & green points.

Hence, regularisation is important to curb w , i.e. make variance smaller, bias \uparrow .

LR classifier : Fairly intuitive to interpret linear decision boundaries.

Not easy to do with non-linear classification boundaries (LR + feature transformation)

Desired :

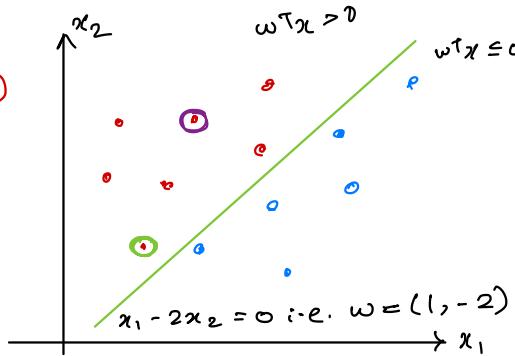
1. The ability to learn complex decision boundaries

2. Interpretability of the model, i.e., how useful are each of the feature/attributes

Decision Tree (DT) classifiers

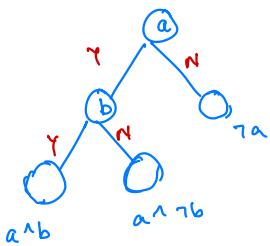
DT is an interpretable model whose final prediction can be written as a "disjunction of conjunctions" based on attribute values over training instances.

- Each path in the tree is a conjunction of attribute based constraints
- Tree : disjunction of the conjunctions



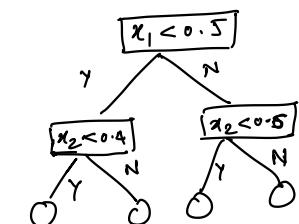
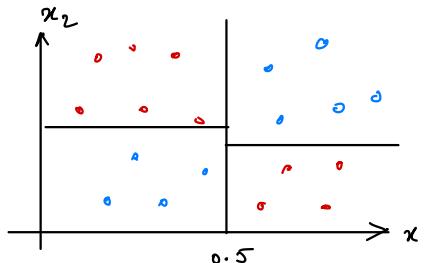
linearly separable data

1) $w = (10^5, -2 \times 10^5)$
is an equivalent solution



$$\text{Tree} \equiv (a \wedge b) \vee (a \wedge \neg b) \vee (\neg a)$$

Decision boundaries of decision trees



DT divides feature space into hyper-rectangles

\equiv the resulting boundaries are axis parallel hyperplanes

Typically, we work with single variable decisions. If we work with linear queries, we get general hyperplanes.

Finding the smallest optimal DT, i.e., optimal w.r.t. some metric involving all attributes \Rightarrow NP-hard.

DT estimation is largely greedy \Rightarrow recursively build the tree.

SIMPLE DT TEMPLATE

1. Start from an empty instance with all instances
2. Pick the best attribute to split on
3. Repeat step 2 recursively on each new node until a stopping criterion is met.

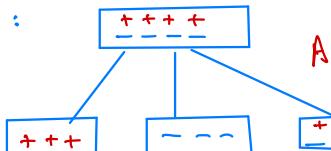
Two important questions

Q1 : What is the notion of a "best" attribute and how to find it ?

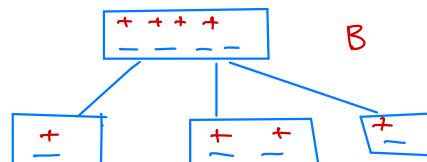
Q2 : what is a good stopping criterion ?

Tree shouldn't be too deep to avoid overfitting

Q1 :



A

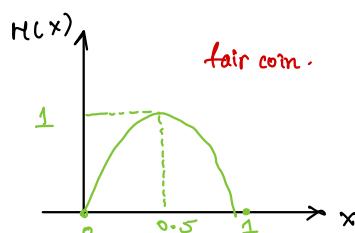


B

Intuition : A good split for an attribute results in subsets that are (nearly) all +ve or all -ve .

Information Gain

$$H(X) = - \sum_{x \in X} P(X=x) \log \underbrace{P(X=x)}_{\text{default base 2}}$$



High entropy \Rightarrow nearly uniform distribution
 \Rightarrow high uncertainty

Low entropy \Rightarrow distribution has well-defined modes
 \Rightarrow low uncertainty

Entropy of a dataset S

$$H(S) = - \sum_{\substack{i \\ \text{set of} \\ \text{attributes}}} p_{i,S} \log p_{i,S}, \text{ where } p_{i,S} \text{ is the relative count of instances in } S \text{ with label } i$$

$$\text{Gain}(S, a) = H(S) - \sum_{v \in \text{values}(a)} \frac{|S_v|}{|S|} H(S_v)$$

where S_v is the subset S whose instances all have attribute a taking the value v .

Lecture

Information Gain of a feature X

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

$$H(Y|X) = \sum_x p(x) H(Y|X=x)$$

$$x^* = \arg \max_x IG(X) = \arg \min_x H(Y|X=x)$$

Few more measures of separation:

- Gini index = $1 - \sum p_i^2$, p_i = fraction of label i in dataset [total = weighted sum]
- Just entropy of the label.

When do we stop?

1. Stop when the data is pure - all elements in dataset are of the same label
2. Stop if number of instances at a node is less than some fixed k
3. Stop if maximum gain at a node \leq threshold. (Used in conjunction with stop if the total number of leaves \geq threshold)

Alternative to stopping criterion : pruning

Reduced error pruning = build tree till end and prune nodes that do not improve error on validation set (prune = remove subtree rooted at node)
{stop down search}

Bootstrap aggregation or Bagging

Consider a dataset D of size n . Create M datasets of same size n from D by sampling with replacement. Bootstrapping - reduces variance. Now, train M models on these datasets. The final prediction is the majority vote over the M models. This is called bootstrap aggregation.

Random Forests

Special case of bagging - they are bagged decision trees. For each DT we also randomly select a subset of features to split the DTs on. This is done to reduce correlation between models - if there was a dominant feature, then all models would separate on this, and there would not be much point in many models. Irrelevant features in decision trees? The majority vote will take care of that.

04/09/23

Lecture

Perceptron Classifier

Linear models of classification. Goal is to estimate hyperplanes that best separate the training data.

$$\text{Eqn of hyperplane} \Rightarrow w^T x = 0$$

Perceptron classifier at test time predicts a test instance x to have the label $h(x) = \text{sign}(w^T x)$

$$\text{where } \text{sign}(z) = \begin{cases} 1, & \text{if } z \geq 0 \\ -1, & \text{if } z < 0 \end{cases}$$

Perceptron is :

1. A basic unit of the neural model of learning [Perceptron is very loosely based on a neuron]
2. Online : It makes updates by seeing training instances one at a time
3. Error-driven or mistake driven : weight updates are triggered only when the perceptron makes a mistake, does nothing otherwise

Perceptron algorithm

Inputs: training instances

Initialise weight vector as w^D , iteration counter $t=0$

for a fixed # of iterations, i.e., $t=0..T-1$

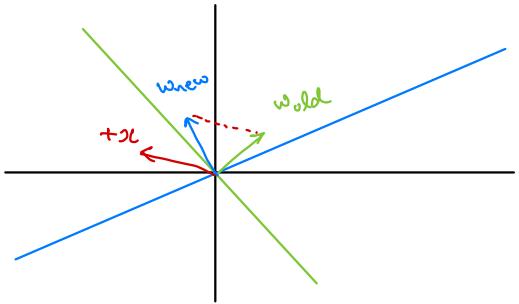
Pick a random training example (x^i, y^i) [or from a shuffled train set sequentially]

$$\text{Predict } \hat{y} = \text{sign}(w_t^T x^i)$$

$$\text{if } \hat{y} \neq y$$

$$w_{t+1} \leftarrow w_t + y^i x^i$$

return w_T



Is the perceptron algorithm guaranteed to improve after a weight update on an erroneous training sample

- Yes, but it is not guaranteed to become correct immediately after an update

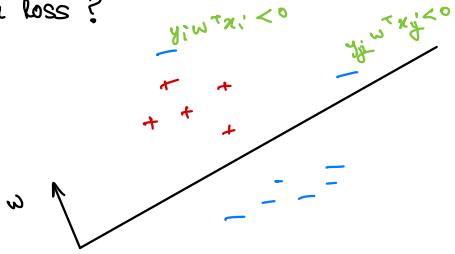
$$y w_{\text{new}}^T x = y (w_{\text{old}} + y \xi)^T x \\ = y w_{\text{old}}^T x + y^2 \|\xi\|_2^2 \\ > y w_{\text{old}}^T x$$

we know
 $y w_{\text{old}}^T x < 0$
and $y w_{\text{new}}^T x > y w_{\text{old}}^T x$

thus making it closer to a positive value.

Since sign is encoded in $y \in \{\pm 1\}$

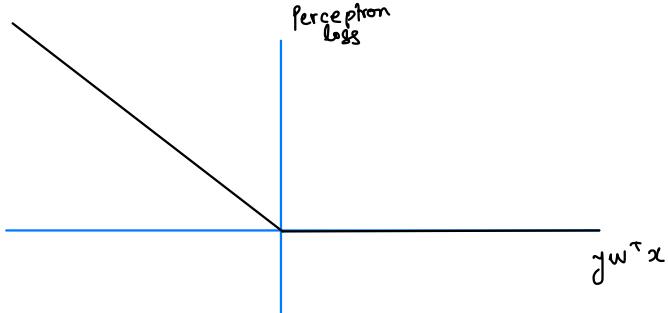
What is the perceptron algorithm minimizing as a loss?



Minimize the number of misclassified examples

$$\Rightarrow \min \sum_{i \in \{\text{indices of misclassified ex.}\}} -y_i w^T x_i$$

Perception loss of $(x, y) = \max \{0, -y w^T x\}$



$$L_{\text{perception}}(w, D) = \sum_i \max \{0, -y_i w^T x_i\}$$

Run SGD on perception loss to find w . For this compute subgradients.
Run SGD on perception loss to find w . For this compute subgradients.
A subgradient of a convex $f^m f$ at w_0 is all vectors g s.t. for any other point w , we have $f(w) - f(w_0) \geq g^T (w - w_0)$

Consider a linearly separable dataset D , i.e. one such that there exists a (unit) vector u such that $y = \text{sgn}(u^T x)$, for every $x, y \in D$.

Without loss of generality, $\forall x_i \in D, \|x_i\| \leq 1$.

Define margin of separation D w.r.t. u be defined as

$$\gamma = \min_{x \in D} |u^T x|$$

that is, γ is the minimum distance of pt $x \in D$ from hyperplane $u^T x = 0$.

Thm: If $\exists u$ s.t. $(x, y) \in D$, we have $y u^T x \geq \gamma$. Then, perceptron algorithm will make no more than $\frac{1}{\gamma^2}$ mistakes on the training dataset when initialized with $w=0$.

Proof: w appears to converge to a loss vector

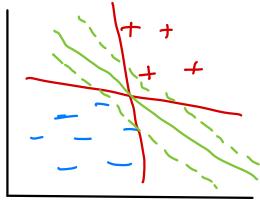
$$\begin{aligned} w_{i+1} &= w_i + y_i x \\ w_{i+1}^T u &= w_i^T u + y_i x^T u \geq w_i^T u + \gamma \\ \|w_{i+1}\|^2 &= \|w_i\|^2 + 2y_i w_i^T x + \|x\|^2 \leq \|w_i\|^2 + 0 + 1 = \|w_i\|^2 + 1 \\ \Rightarrow \|w_k\| &\leq \sqrt{k} \quad \forall k. \end{aligned}$$

If perceptron algorithm makes k mistakes,

$$\sqrt{k} \geq \|w_k\| \geq w_k^T u \geq k\gamma \Rightarrow k \leq \frac{1}{\gamma^2}$$

$$\begin{aligned} \underline{\text{M2}}: \quad y w_{i+1}^T x &= y w_i^T x + \|x\|^2 \leq y w_i^T x + 1 \\ \max(0, -y w_{i+1}^T x) &\leq \max(0, -y w_i^T x) - 1 \end{aligned}$$

Support Vector Machines



SVMs are max-margin classifier that aim to maximize the margin while correctly classifying all the training points. Here, margin is the distance of the closest point across both classifiers from the hyperplane.

$$\text{Margin } \gamma(w, b) = \min_i \frac{|w^T x_i + b|}{\|w\|} \rightarrow \textcircled{A}$$

SVM wants to maximize the margin & correctly classify all training pts

$$\max_{w, b} \underbrace{2 \gamma(w, b)}_{\text{for a nice gradient}} \text{ s.t. } \forall i \quad y_i (w^T x_i + b) \geq 0 \rightarrow \textcircled{B}$$

Substituting \textcircled{A} in \textcircled{B} ,

$$\max_{w, b} \underbrace{\frac{2}{\|w\|}}_{\text{from } \textcircled{C}} \min_i (w^T x_i + b) \text{ s.t. } y_i (w^T x_i + b) \geq 0 \rightarrow \textcircled{C}$$

Fix the scale of w, b s.t. $\min_i (w^T x_i + b) = 1 \rightarrow \textcircled{D}$

Note: This is allowed as $d(x, L) = \frac{\|L(x)\|}{\sqrt{\text{sum of coeff}^2}}$

from \textcircled{C} and \textcircled{D} ,

$$\max_{w, b} \frac{2}{\|w\|} \text{ s.t. } y_i (w^T x_i + b) \geq 0, \quad \min_i (w^T x_i + b) = 1$$

$$\equiv \min_{w, b} \frac{1}{2} \|w\|^2 \text{ s.t. } y_i (w^T x_i + b) \geq 0, \quad \min_i (w^T x_i + b) = 1 \rightarrow \textcircled{E}$$

Hard margin sum optimization problem

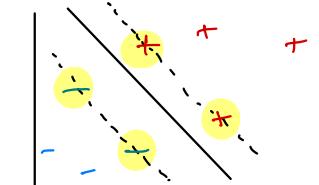
from \textcircled{E} ,

$$\begin{aligned} & \min_{w, b} \frac{1}{2} \|w\|^2 \\ & \text{s.t. } y_i (w^T x_i + b) \geq 1 \end{aligned}$$

Note: equivalent to \textcircled{E} ! (why?)

$$\begin{aligned} y_i \in \{-1, 1\} \text{ Hence, } y_i (w^T b + 1) &\geq 1 \\ \Rightarrow (w^T b + 1) &\geq 1 \\ \Rightarrow \min |w^T b + 1| &\geq 1. \end{aligned}$$

Why hard margin? Because we want every training example to satisfy margin constraint



Highlighted points = tight constraints. If not, we can rescale w, b s.t. constraints become tight. The tight points are "support vectors".

Nob Linearly Separable data

⇒ This means that ∃ some pts which variable and margin constraints

$$\exists i \quad y_i(w^T x_i + b) < 1 \quad \{ \text{given non-lin sep. data for any } w \}$$

$$\Rightarrow \exists i \quad (1 - y_i(w^T x_i + b)) > 0$$

$$\text{Constraint violation} = \begin{cases} 0 & \text{if } 1 - y_i(w^T x_i + b) \leq 0 \\ 1 - y_i(w^T x_i + b) & \text{otherwise.} \end{cases}$$

minimize the total violation of margin across D_i :

$$\sum_i \max\{0, 1 - y_i(w^T x_i + b)\} \quad (\text{Hinge loss})$$

revised unconstrained SVM optimization problem

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_i \max\{0, 1 - y_i(w^T x_i + b)\} \quad (\text{Regularised Hinge loss})$$

Soft-margin sum optimization problem

Equivalent constrained opt. to A

$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_i \xi_i \rightarrow \text{slack variables.}$$

$$\text{s.t. } \forall i \quad y_i(w^T x_i + b) \geq 1 - \xi_i;$$

$$\forall i \quad \xi_i \geq 0$$

What is the role?

C is a penalty or regularization factor. Balances the trade-off between large margins and large-training error.

SVM Optimization Problem

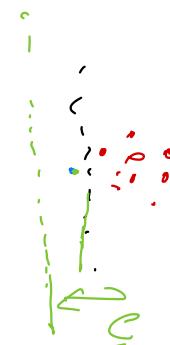
$$\min_{w, b} \frac{1}{2} \|w\|^2$$

$$\text{s.t. } y_i(w^T x_i + b) \geq 1$$

$$\text{Define } g_i(w) = 1 - y_i(w^T x_i + b) \leq 0$$

For each constraint, define a Lagrange multiplier $\alpha_i \geq 0$, write Lagrangians.

$$L(w, \alpha) = \frac{1}{2} \|w\|^2 + \sum_i \alpha_i (1 - y_i(w^T x_i + b))$$



generic constrained optimization problem is :

$$\min_w f(w) \quad (\text{Primal})$$

$$\text{s.t. } g_1(w) \leq 0, \dots, g_n(w) \leq 0$$

$$\text{Lagrangian } L(w, \alpha) = f(w) + \sum \alpha_i g_i(w)$$

Lagrange mult

equivalent formulation of the primal is

$$\min_w \max_{\alpha \geq 0} L(w, \alpha) \quad [\text{removing constraints on } w]$$

why the equivalence ?

$$\max_{\alpha \geq 0} f(w) + \sum \alpha_i g_i(w) = \begin{cases} f(w) & \text{if } w \text{ satisfies} \\ & \text{all constraints} \\ \infty & \text{otherwise} \end{cases}$$

Consider the dual problem

$$d^* = \max_{\alpha \geq 0} \min_w L(w, \alpha)$$

In general $d^* \leq p^*$.

$$\boxed{d^* = p^*}$$

Dual for SVM

$$\max_{\alpha \geq 0} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

$$\text{s.t. } \forall i \alpha_i \geq 0, \sum_i \alpha_i y_i = 0$$

(Hard-margin SVM)

In soft-margin SVM,

$$0 \leq \alpha_i \leq C.$$

Everything else is same as hard margin

$$L(w, \alpha) = \frac{1}{2} \|w\|^2 + \sum_i \alpha_i (1 - y_i (w^T x_i + b))$$

$$\nabla_w L(w, \alpha) = w + \sum_i \alpha_i (-y_i x_i) = 0$$

$$\Rightarrow w = \sum_{i=1}^n \alpha_i y_i x_i$$

Similarly

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{i=1}^n \alpha_i y_i = 0$$

put in
 $L(w, \alpha)$

Observations based on the dual form,

① $w = \sum_{i=1}^n \alpha_i y_i x_i$ $\leftarrow \alpha_i > 0 \Rightarrow$ tight constraints, i.e. sv , hence w depends only on support vectors.

② If the solution to primal is w^* and to dual is some α^* then KKT conditions on w^* and α^* hold for both SVM problems.

One of the dual complementarity constraint, is of interest

$$\alpha_i^* y_i g_i(w^*) = 0 \quad \forall i$$

$\alpha_i > 0$ corresponds to those training examples for which $g_i(w)$ is the equality, i.e. $g_i(w) = 0$.

$$\Rightarrow y_i (w^T x_i + b) = 1$$

These are the "support vectors"

Dual Hard-margin SVM

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i^T x_j, \quad \alpha \geq 0, \quad \sum_i \alpha_i y_i = 0$$

③ The dual form operates only on inner products of training examples pairs.

Test Time

$$\begin{aligned} \text{Prediction} &= \text{sgn}(w^* T x + b^*) \\ &= \text{sgn}\left(\left(\sum_i \alpha_i^* y_i x_i\right)^T x + b^*\right) \\ &= \text{sgn}\left(\sum_i \alpha_i^* y_i \underbrace{x_i^T x}_{\text{again, inner product}} + b^*\right) \end{aligned}$$

If we want a non-linear decision boundary using SVMs,

$$\text{prediction} = \text{sign}\left(\sum_i \alpha_i^* y_i \underbrace{\phi(x_i)^T \phi(x)}_{\text{feature transform } \phi: \mathbb{R}^d \rightarrow \mathbb{R}^m} + b^*\right)$$

Two problems :

1. How do I find $\phi(x)$?

2. Computationally expensive to calculate $\phi(x)$ and compute dot product

Kernels

Consider a kernel function $K(x, y) = \phi(x)^T \phi(y)$

If we can compute $K(x, y)$ without computing $\phi(x)$, then this is useful (kernel trick)

Consider the following example, for instances $x = (x_1, x_2), y = (y_1, y_2)$

$$\phi(x) = \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \end{bmatrix}.$$

Define a kernel function $K(x, y) = (x^T y)^2$

Verify that $K(x, y) = \phi(x)^T \phi(y)$

Can we use any function for K ? No

① A kernel $K(x, y)$ is valid kernel function if it corresponds to an underlying $\phi(x)$

[projecting x to a

that is, $K(x, y) = \phi(x)^T \phi(y)$

② Let $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ is a valid kernel, the necessary and sufficient condition is that for any finite set $\{x_1, \dots, x_n\}$ the kernel matrix $[K(x_i, x_j)]$ is symmetric and positive semi-definite [Mercer's theorem]

③ Use properties of kernels over existing valid kernels

$$K(x, y) = K_1(x, y) + K_2(x, y)$$

$$K(x, y) = K_1(x, y) K_2(x, y)$$

$$K(x, y) = \alpha K_1(x, y), \alpha > 0$$

Examples of valid kernels

1. Polynomial kernel $K(x, y) = (x^T y)^d$ [or $(c + x^T y)^d$: degree $\leq d$]

2. Radial basis fn / Gaussian kernel : $K(x, y) = e^{-\frac{\|x-y\|^2}{\sigma^2}}$

measure of
similarity b/w
pairs of points

Q : Is $K(x, y) = \|x+y\|^2$ a valid kernel?

$$K(x, y) = (x+y)^T (x+y)$$

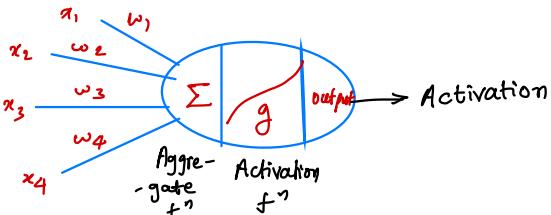
$$K(\omega, 0) = 0 \Rightarrow \phi(0) = 0.$$

$\Rightarrow K(x, 0) = 0$, but $K(x, y) = x^T y$. (contradiction)

12/09/23

Lecture

McCulloch-Pitts Neuron Model

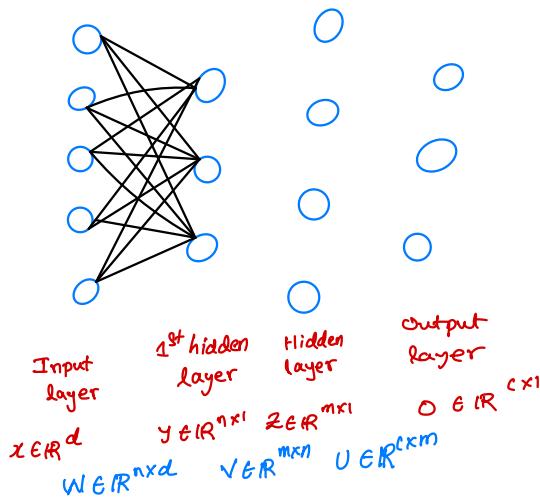


limitation: classification of xor

$$\text{Neuron activation} = g\left(\sum_{i=1}^n w_i x_i\right)$$

$$\Sigma = \sum, g = \text{sign} \rightarrow \text{Perceptron}.$$

Let's stack neuron models together to form a multi-layered neural network



w is an $n \times d$ matrix where w_{ij} refers to weight of connection from i^{th} row in input layer to j^{th} row in hidden layer

$$y = g(wx + b), b = \text{bias setting}$$

↑
applied
elementwise

$$z = g(y + b') = g(\gamma g(wx + b) + b')$$

$$o = g(uz + b'')$$

Fully-connected (or) feed-forward network
(not necessarily fully connected, but this is default)

Output $o \in ℝ^{cx1}$: What are useful loss functions?

① Regression : Squared error loss = $\sum_i \frac{1}{2} (o_i - y_i)^2$

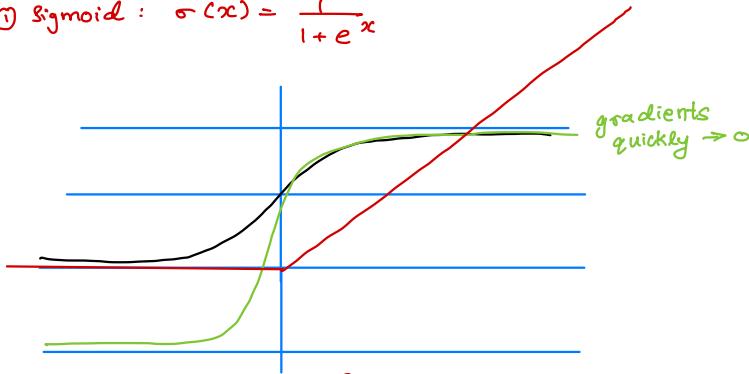
② Classification : A very handy transformation is the SOFTMAX operation.

$$\text{SOFTMAX}(o_j) = \frac{e^{o_j}}{\sum_{k=1}^c e^{o_k}}, j \in \{1, \dots, c\} \quad \{ \text{generalisation of sigmoid} \}$$

$p_{NN}(y|x) = \text{SOFTMAX}(o)$ {then we cross-entropy loss} on this

Common Activation Functions (g)

① Sigmoid: $\sigma(x) = \frac{1}{1+e^{-x}}$



② tanh: $\tanh(x) = \frac{e^{2x}-1}{e^{2x}+1}$ outputs b/w -1 and +1

③ RELU: $\text{relu}(x) = \max(0, x)$ Does remarkably well.
Simple gradients, converges faster

Loss function $f(w, D)$

(credits : KNA)

Typical loss function

For regression problems (the output vector is scalar), with true output y ,
output of MLP (multi-layer perceptron) \hat{y} ,

$$L(y, \hat{y}) = \frac{1}{2} (y - \hat{y})^2 \quad (\text{squared error loss})$$

For classification problems, cross-entropy loss is typically used.

Suppose we have k -classes, output vector of MLP is k -dimensional.

Let y be the true label of the input x , let \hat{y} be output of MLP on input x

$$L(y, \hat{y}) = - \sum_{k=1}^K y_k \log \hat{y}_k \quad (\text{cross-entropy loss})$$

A trick to convert output as probabilities

$$f(y=i|x) = \frac{e^{o_i}}{\sum_{j=1}^k e^{o_j}}$$

where (o_1, \dots, o_K) is the output of the network on input x .

Backpropagation

To obtain weights as function of training data, we use regular SGD.

Initialize weights and biases randomly

while not converged do

 Sample a minibatch of data D

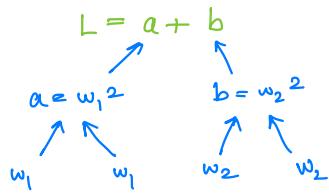
 Compute the gradient of the loss function $\nabla_w L$ on D with respect to weights

 Update the weights, biases using gradient $w \leftarrow w - \alpha \nabla_w L$

end while

Computation of gradients? Use backpropagation.

We can write dependence of loss on weights as directed acyclic graph of loss on the weights e.g. $L(w_1, w_2) = w_1^2 + w_2^2$



$$\text{Here, } g(a, b) = a + b, \alpha = w_1^2$$

$$\begin{aligned} \text{we get } \frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial a} \frac{\partial a}{\partial w_1} \\ &= 1 \cdot 2w_1 \end{aligned}$$

$$\text{similarly } \frac{\partial L}{\partial w_2} = 2w_2$$

Thm (Multivariable Chain Rule)

Suppose $f = f(y_1, \dots, y_n)$ is a function. For some i , suppose $y = g(x_i)$

$$\text{Then, } \frac{\partial f}{\partial x_i} = \sum_{j=1}^n \frac{\partial f}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

So, the idea is we can compute $\frac{\partial L}{\partial w}$ if we know $\frac{\partial u}{\partial w}$ for every child u ,

and form a recursive procedure.

Algorithm

Algorithm has a forward pass, backward pass. Forward pass computes values of nodes in the DAG. The backward pass computes derivatives of the loss with respect to weights (i.e. gradients)

Forward pass Compute the values of nodes in DAG in topological order

Algorithm 2 forward pass

- 1: The values of the leaves are the input values and the values of the weights.
- 2: **for** each node u in the DAG in topological order **do**
- 3: compute $u.val = f(u.child_1.val, u.child_2.val, \dots, u.child_k.val)$
- 4: **end for**

Backward pass

Here, we use the computed values to compute gradient of every node with respect to the leaves (i.e. we don't care abt $\frac{\partial L}{\partial X}$)

While building the DAG, we store dependencies of the parent $di = \frac{\partial u}{\partial (u \cdot \text{child}_i)}$ on its children's functions. We use some of input vals (e.g. $w^T x$), but we do not include them in children because they aren't vars we're differentiating with respect to.

Algorithm 3 backward pass

- 1: The function dw for the leaves (weights) is set to 1.
 - 2: **for** each node u in the DAG in topological order **do**
 - 3: compute $u.\text{grad} = \sum_{i=1}^k u.di \times \text{child}_i.\text{grad}$ (where $.\text{grad}$ is a vector with dimension equal to the number of weights)
 - 4: **end for**
-

Regularization

1. L2 regularization : Add term $\lambda \|w\|_2^2$ to loss

2. Dropout :

During training, randomly set the activations of some neurons to 0. This prevents network from relying too much on any one neuron, and thus prevents overfitting.

- Do not drop out input, output neurons
- At test time, multiply activations of each neuron by probability that it was dropped out during training, i.e. expected value of activation is used.

Remark : An extreme version of dropout is to randomly drop out entire layers of neurons (called drop connect). This makes use of residual connections, which connect the input of a layer to the output of a layer that is not the immediate previous layer.

3. Early stopping :

Stop training when the validation loss starts increasing. This is a form of regularisation because it prevents network from overfitting to training data (we can stop training before it can do so).

Optimizers

1. Momentum

- Large weight updates in complex loss landscapes makes training unstable
- At flat points like saddle points, SGD gets stuck

SGD with momentum addresses these issues.

Let $0 \leq \beta < 1$ and define velocity v_t as

$$v_t = \eta (\nabla_w L(w_t)) + \beta v_{t-1}$$

$\nabla_w L(w_t)$ means $\nabla_w L$ at $w=w_t$, it's not multiplication :)

$$\begin{aligned} \text{so, } v_t &= \eta (\nabla_w L(w_t)) + \beta (\eta \nabla_w L(w_{t-1})) + \beta^2 v_{t-2} \\ &= \dots \\ &= \eta \sum_{i=0}^t \beta^i \nabla_w L(w_{t-i}) \end{aligned}$$

and weight update is,

$$w_{t+1} = w_t - v_t$$

Note: This escapes flat points, but will escape minima as well, then come back as velocity reduces, velocity will likely be small upon coming back so won't overshoot again.

2. ADAGrad (adaptive gradient)

- Have different learning rate for each parameter as some parameters are more sensitive than others

\odot = element-wise multiplication

$$s_t = s_{t-1} + g_t \odot g_t, g_t = \nabla_w L(w_t)$$

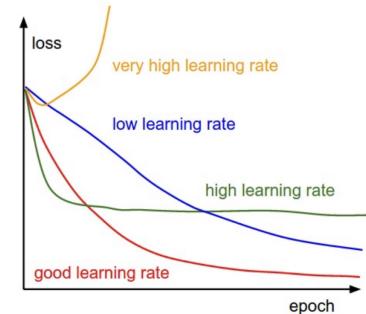
and, weight update is,

$$w_{t+1} = w_t - \eta \frac{1}{\sqrt{s_t} + \epsilon} \odot g_t$$

so, we get

$$\eta_i = \frac{\eta}{\sqrt{\sum_{t=1}^T (\nabla_w L(w_t))^2} + \epsilon}$$

, weight update $(w_{t+1})_i = (w_t)_i - \eta_i (\nabla_w L(w_t))_i$



Issue: learning rate has sum of squares in den^r, so algorithm eventually stops learning.
This issue is solved by the next optimizer

3. RMS Prop

$$s_t = r s_{t-1} + (1-r) g_t \odot g_t$$

Note : This allows us to forget old gradients. This was proposed in a Coursera course by Geoff Hinton !

4. ADAM

Here, we combine the above ideas.

$$v_t = \beta v_{t-1} + (1-\beta) g_t$$

$$s_t = r s_{t-1} + (1-r) s_t \odot s_t, \text{ where } g_t = \nabla_w L(w_t)$$

$$\text{update step: } w_{t+1} = w_t - \frac{\eta}{\sqrt{s_t + \epsilon}} \odot \hat{v}_t$$

$$\text{where } \hat{s}_t = \frac{s_t}{1-r^t}, \quad \hat{v}_t = \frac{v_t}{1-\beta^t}$$

} to correct the fact that, v_t, s_t initialised to 0, so they are biased toward zero.

We have :

$$v_t = (1-\beta) \sum_{i=0}^{t-1} \beta^i g_{t-i}, \text{ with sum of coeff being } 1-\beta^t,$$

so it helps to make velocity directly a moving avg of gradients for small t .

Optimizers

- A. SGD with momentum : smoothed fast gradients over time
- B. ADAGrad : Adaptive learning rates per feature
- C. RMS Prop
- D. Adam : Combination of A and B

$$\text{ADAM} \quad s_t \leftarrow r s_{t-1} + (1-r) g_t \odot g_t$$

$$v_t \leftarrow \beta v_{t-1} + (1-\beta) g_t$$

$$\hat{s}_t = \frac{s_t}{1-r^t}, \quad \hat{v}_t = \frac{v_t}{1-\beta^t} \quad \rightarrow \text{Bias correction}$$

$$\text{wt update for Adam} \quad w_t \leftarrow w_{t-1} - \frac{1}{\sqrt{s_t} + \epsilon} \odot \hat{g}_t$$

$$s_t = (1-r) g_t \odot g_t + r s_{t-1}$$

$$= (1-r) g_t \odot g_t + (1-r) r g_{t-1} \odot g_{t-1} + \dots + (1-r) r^t s_0$$

Sum of coeff of s_t : $(1-r)(1+r+\dots) \approx 1$ (if t is large)

for small t , sum = $1 - r^t$

So, $\frac{s_t}{1-r^t}$ mitigates the defect.

Consider the problem of recognising an object in an image

DESIDERATA

- (A) Locality
- (B) Translation invariance
- (C) Parameter efficient

FFN

X (fully connected)

CNNs



X (different connections)

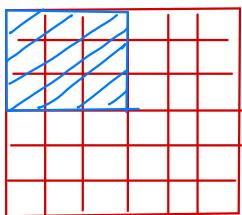


X (expensive)



FFN = feed forward NN

Main operation underlying CNNs is convolution.



(A)



Window multiplication = "cross-correlation"
but not convolution.

For this purpose, we loosely use the term convolution.

(B)

$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$
$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$
$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$

Mean Kernel

= Blurred image

Monochromatic
+
no change
Variety
↓
"average"

B

$$\begin{matrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{matrix} = \begin{array}{l} \text{Detects} \\ \text{vertical} \\ \text{edges} \end{array}$$

Monochromatic = zero

Sobel filter = edge-detection filter

CNNs let you learn the kernel weights from data

Ref: Conv slides

Batch Normalisation

Recall that preprocessing /standardizing the inputs to have zero mean & unit variance was a typical first step. Why?

① Constrains the function complexity of the hypothesis in ge

② Standardization works well with GD-style optimizers by putting parameters a priori in a smaller space.

In deep NN can we add normalization within network

Batch Normalization ('15)

Variables can have val with largely varying magnitude .

BN authors conjectured that such a drift in distribution of internal variables (dubbed "internal covariate shift") can cause convergence challenges in deep networks

neurips '18: just ensures smoothness , cleaner updates (no covariate shift)
(worse covariate shift, smoother \Rightarrow better acc %.)

Adaptive solvers (Adagrad) aim to address this issue . Alternate
is to pre-emptively prevent this from happening (e.g. Batch Norm)

Batch Norm has 2 steps

1. In each training iteration, from the activations, subtract mean and divide by std-deviation. Both mean, σ computed using batch level statistics
2. Next, apply scaling coeff, shift by offset (learnable)

Let B denote the batch, $x \in B$ then

$BN(x)$ can be written as

$$BN(x) = \gamma \odot \frac{x - \hat{\mu}_B}{\hat{\sigma}_B} + \beta \quad \text{learnable}$$

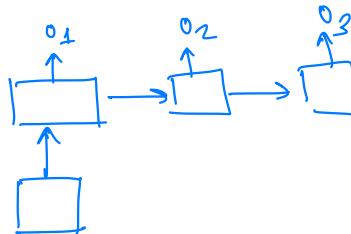
\uparrow
elementwise product

$$\hat{\mu}_B = \frac{\sum_{x \in B} x}{|B|}, \quad \hat{\sigma}_B = \sqrt{\frac{1}{|B|} \sum_{x \in B} (x - \hat{\mu}_B)^2} + \epsilon$$

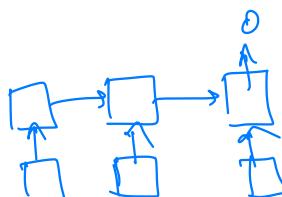
At test time : Use running avg of mean, s.d.-values during training

Note : BN is applied after affine, before non-lin. activation

(BN is linear, but causes "smoothness" in traversal of loss landscape)

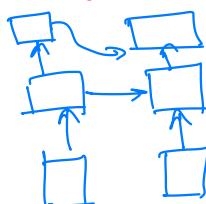


fixed input
variable length output
e.g. image descriptor

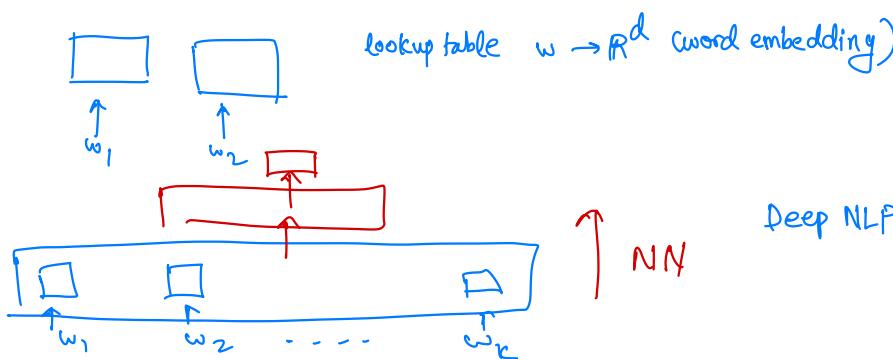


fixed o
variable length input
e.g. sentiment analysis

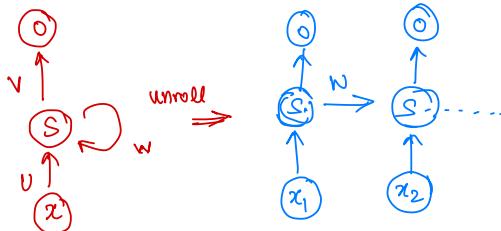
Both variable length (e.g. sequence labelling problem) . output length = $f(\text{input len})$



Language Modelling (causal language problem)
Given a word history w_1, \dots, w_{t-1} what is $w^* = \arg \max_w P(w_t | w_1, \dots, w_{t-1})$?



Recurrent NNs



$$s_t = \tanh(u s_{t-1} + v x_t + b)$$

$$o_t = V s_t + b$$

$$y_t = \text{softmax}(o_t)$$

RNNs by design model temporal dependencies

Challenge: Vanishing and exploding gradients.

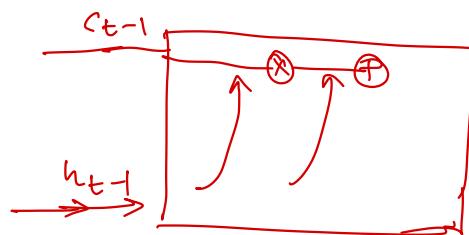
Recall that gradients in earlier layers = prod of grad in later layers (exploding gradients)

Band-aid: clip gradients \leq threshold

If latter takes small vals: vanishing gradients

Long short term memory (LSTM) network

In LSTMs, the recurrent unit is enhanced with a cell state (or a memory state) and a number of learnable gates



forget gate¹

$$f_t = \sigma(w_f[x_t; h_{t-1}] + b_f)$$

(element wise)

$\rightarrow 0$: forget

$\rightarrow 1$: retain

$$\text{input gate: } i_t = \sigma(w_i[x_t, h_{t-1}] + b_i)$$

$$\text{modified state: } \tilde{c}_t = \tanh(w_c[x_t, h_{t-1}] + b_c)$$

Final state: c_t

$$c_t = f_t \odot c_{t-1} + i_t \odot \tilde{c}_t$$

$$\text{output gate: } o_t = \sigma(w_o[c_t, h_{t-1}] + b_o)$$

$$\text{hidden state: } h_t = o_t \odot \tanh(c_t)$$

GRU: gated recurrence unit.

clustering

Consider a set of points $x_1, \dots, x_n \in \mathbb{R}^d$. Assume that each datapoint x_i belongs to one of K clusters (we want points within each cluster to be "similar", i.e. close to one another in terms of Euclidean distances)

How do we identify such clusters of datapoints?

E.g. Stock market indices, demographics, social networks

Motivation:

- (A) Prediction based on labelling single point
- (B) Outlier detection
- (C) (Lossy) compression

k-means clustering

Main intuition: Assume that there are K clusters. every datapoint is close to the centre or the mean of its assigned cluster.

If we fix the cluster centres, it's easy to make assignments for every datapoint.

If we fix assignments, easy to compute cluster centres.

(chicken & egg problem)

Fix one, optimize for the other and alternate.

Fix one, optimize for the other and alternate.

Note: This is an "alternating minimization algorithm".

k-means objective

Find cluster centres $u_1, \dots, u_K \in \mathbb{R}^d$ and assignments c^1, \dots, c^n [where each c_i is a one hot encoding of size K] such that sum of squared distances of all datapoints $x_1, \dots, x_n \in \mathbb{R}^d$ from their assigned cluster centres is minimized

$$\text{formally: } \sum_{i=1}^n \sum_{k=1}^K c_k^i \|u_k - x_i\|^2$$

where $c_k^i = \underline{1} [x_i \text{ assigned to cluster } k]$
 "best assignment"

k-means Heuristic : Fix one, optimize the other

Step 1 : Fix μ , optimize the assignment.

for the i th datapoint, $\sum_{k=1}^K c_k^i \|\mu_k - x_i\|^2$

$$c_k^i = \begin{cases} 1, & \text{if } k = \operatorname{argmin}_j \|\mu_j - x_i\|^2 \\ 0, & \text{otherwise} \end{cases}$$

Step 2 : Fix C , optimize cluster centres

For $k=1, \dots, K$

$$\frac{\partial}{\partial \mu_k} \sum_i \sum_k c_k^i \|\mu_k - x_i\|^2 = 0$$

$$\Rightarrow \mu_k = \frac{\sum_i c_k^i x_i}{\sum_i c_k^i} \quad \text{*e.g. centroid of cluster points*}$$

K-means Algorithm

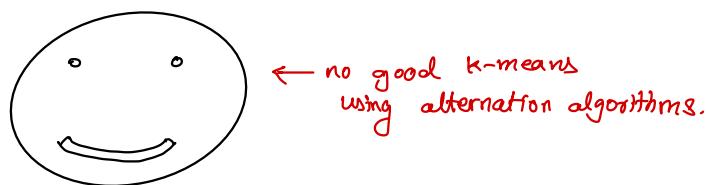
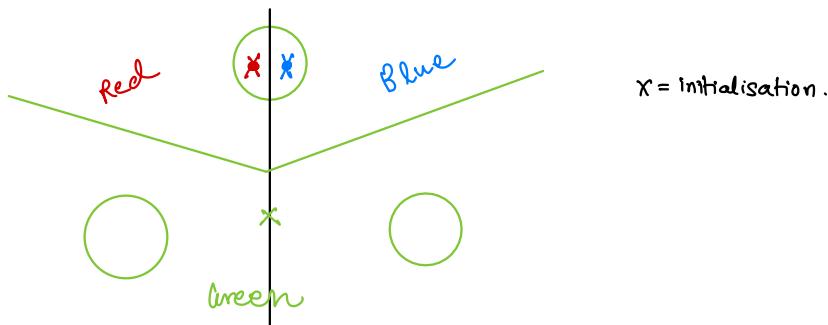
Initialize the cluster centres randomly (μ_1, \dots, μ_K)

Repeat until convergence [i.e. assignments do not change]

Assignment : Assign each datapoint to nearest cluster

Update : Recompute the cluster centres as the mean of all datapoints assigned to a cluster.

Note : This is not guaranteed to converge "correctly"



Note : K-means is NP-hard.

Q. Is K-means guaranteed to converge? terminate.

K-means reduces the cost at every iteration (update as well as assignment, both local optima)
this, combined with the fact that number of cluster assignments is finite (since finite set of pts).

Hence, cycles cannot happen in configurations.

Convergence (termination) is guaranteed.

So far,

K-means is a "HARD" assignment meaning each data point is assigned to a single cluster

soft assignment? (probabilistic)

$$c_k^i = \frac{\exp(-\beta \| \mu_k - \mathbf{m}_i \|^2)}{\sum_j \exp(-\beta \| \mu_j - \mathbf{m}_i \|^2)}$$

Note: $\beta \rightarrow \infty$ then this is the hard k-means problem.

Lecture

10/10/23

Abstract alg.

Θ = assignments

(ref: abstract code)

ϕ = cluster centres

1. Both updates cost decreases, so cost always decreases

* 2. True

3. False. Θ_0 can be revisited. Θ can be repeated almost twice

4. Θ is finite, can be repeated almost twice, so always terminates.

Recap: Hard k-means vs Soft k-means

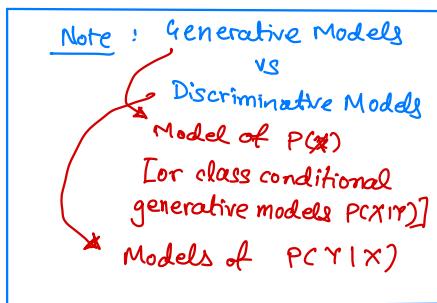
Generative/Probabilistic Model of Clustering

Imagine that dataset $D = \{x_1, \dots, x_N\}$ are sampled from a generative model.

Idea: once you have $P(X)$, we can generate joint distributions.

NNB: discriminative model.

Generative AI: (Andrew Ng: logistic regression & its generative counterpart)



← nice reference.

(add" reading)

Note: You can use $p(x|y)$, $p(y)$ for discriminative models (bayes rule)

Then, estimate the parameters of this generative model using MLE

Clustering using a Generative Model

Consider a datapoint x from the dataset D

- Pick a cluster z from $\{1, \dots, K\}$ s.t. $P(z=k) = \pi_k$
- Given a z , sample from a gaussian distribution $p(x|z) = N(x; \mu_z, \Sigma)$

[Note: For simplicity of subsequent derivations, we'll assume $\Sigma = \mathbb{I}$;
more generally, $\Sigma = \Sigma_z$]

Generative Process :

$$P(z=k) = \pi_k, \quad \sum_{k=1}^K \pi_k = 1$$

$$p(x|z=k) = N(x; \mu_k, \mathbb{I})$$

How do we estimate $\{\pi_k, \mu_k\}_{k=1}^K$?

use MLE (max. likelihood estimation)

$$\max \log p(x) = \log \sum_k P(z=k, x) \quad \begin{matrix} \text{observed vars} \\ \uparrow \text{hidden variable.} \end{matrix}$$

$$\begin{aligned} \text{single datapoint} &= \log \sum_k p(x|z=k) P(z=k) \\ &= \log \sum_k \pi_k N(x; \mu_k, \mathbb{I}) \end{aligned}$$

mixture component weights
mixing coefficients

$$p(x) = \sum_k \pi_k N(x; \mu_k, \mathbb{I})$$

Gaussian mixture model
(GMM)

Aside : Note on mixture models.

Going back to the objective to maximize for generative model

$$\sum_{i=1}^N \log p(x_i) = \sum_{i=1}^N \log \left(\sum_{k=1}^K \pi_k N(x_i; \mu_k, \mathbb{I}) \right)$$

this is cumbersome to optimize.

So, let's simplify to a hard cluster setting, i.e., for every x_i , assume you know corresponding z_i (z_i are also given) \leftarrow Fully observable setting

$$\begin{aligned} \max \sum_i \log p(x_i) &= \sum_i \log \sum_k P(z=k, x_i) \quad \begin{matrix} \# \text{ to do:} \\ \downarrow \text{fill in steps} \end{matrix} \\ &= \max \sum_i \sum_k \mathbf{1}(z_i=k) (\log N(x_i; \mu_k, \mathbb{I}) + \log \pi_k) - \textcircled{A} \end{aligned}$$

Do MLE now for μ_k, π_k for \textcircled{A}

$$\mu_k = \frac{\sum_i \mathbf{1}(z_i=k) x_i}{\sum_i \mathbf{1}(z_i=k)}, \quad \pi_k = \frac{1}{N} \sum_i \mathbf{1}(z_i=k) \quad \begin{matrix} \# \text{ to do:} \\ \downarrow \text{fill in steps} \end{matrix}$$

Now, we build up to a solution where we don't know what z_i 's are.

Consider the objective in ④

We do not have information about $\mathbb{1}[z_i=k]$, but, we can compute the $\mathbb{E}[\mathbb{1}[z_i=k]] = P(z_i=k|x_i) \leftarrow \# \text{ To do: } \mathbb{E} \text{ over what? fill in explanation}$

so, compute $P(z_i=k|x_i) = \frac{P(x_i|z_i=k) P(z_i=k)}{P(x_i)}$ [refer to this as $r_{k,i}$]
 $= \frac{N(x_i; \mu_k, \Sigma)}{\sum_j \pi_j N(x_i; \mu_j, \Sigma)}$

Looking at expression in ④, replace indicator with $r_{k,i}$ to get : Note :

$$\sum_i \sum_k r_{k,i} (\log N(x_i; \mu_k, \Sigma) + \log \pi_k) - \textcircled{B}$$

Get ④ from ④ by assuming $r_{k,i}$ are fixed,

meaning they use current values of parameters μ_k, π_k

chicken-egg situation once again, to get r we need \textcircled{A} and vice-versa.
 $r \Leftrightarrow \pi_k, \mu_k$

This motivates the expectation maximization algorithm [EM algorithm]

Lecture

12/10/23

Recall the fully observable setting, the MLE

$$\sum_i \log P(x_i, z_i) = \sum_i \sum_k \mathbb{1}[z_i=k] (\log \pi_k + \log N(x_i; \mu_k, \Sigma))$$

In the real setting only x_1, \dots, x_N are observed, z_i 's are hidden/latent

$$r_{k,i} \triangleq P(z_i=k|x_i) = \pi_k N(x_i; \mu_k, \Sigma) / Z$$

If we treat $r_{k,i}$ values as fixed, then, MLE objective becomes

$$\sum_i \sum_k r_{k,i} (\log \pi_k + \log N(x_i; \mu_k, \Sigma))$$

$\nabla_{\mu_k} \text{LLg}(\mu_k; \pi_k) = 0$, solve for μ_k, π_k

$$\mu_k = \frac{\sum_i r_{k,i} x_i}{\sum_i r_{k,i}}, \quad \pi_k = \frac{\sum_i r_{k,i}}{N}$$

weight that
 x_i commits
to μ_k

This motivates EM algorithm.

EM iteratively finds a lower bound of the log likelihood and optimizes it.

Initialise: $\mu_k^{cur} = \mu_k^{init}$, $\pi_k^{cur} = \pi_k^{init} \quad \forall k \in \{1, 2, \dots, n\}$

E-step : Compute posterior probas.

$$\pi_{k,i} = \frac{\pi_k^{cur} N(x_i; \mu_k^{cur}, \Sigma)}{\sum_j \pi_j^{cur} N(x_i; \mu_j^{cur}, \Sigma)}$$

M-step: Estimate μ_k, π_k

$$\mu_k^{new} = \frac{\sum_i \pi_{k,i} x_i}{\sum_i \pi_{k,i}}, \quad \pi_k^{new} = \frac{\sum_i \pi_{k,i}}{N}$$

Check if the log-likelihood value is within a threshold of the previous log-likelihood

If yes, converged \rightarrow terminate

else $\mu_k^{curr} = \mu_k^{new}, \pi_k^{curr} = \pi_k^{new}$ repeat E step and M step

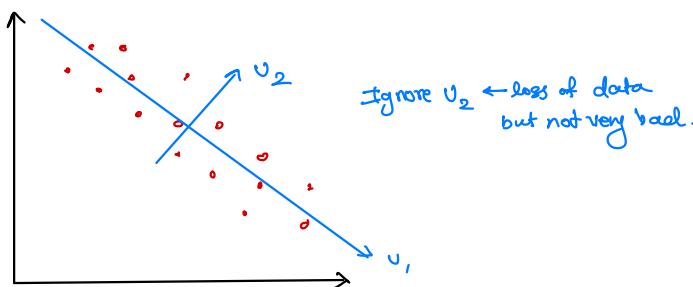
Dimensionality Reduction

Most popular dimension reduction technique is principal component analysis (PCA)

Consider a point $x \in \mathbb{R}^d$. Aim of PCA is to find a low dimensional projection of x , i.e.

$$z = u^T x, \quad u \in \mathbb{R}^{d \times k} \text{ where } k < d.$$

How do we estimate u ?



Minimize reconstruction error

$$\text{Encode: } z = u^T x$$

$$\text{Decode: } \hat{x} = u z = u u^T x$$

Reconstruction Objective fn

$$\min_{U^T U = I} \sum_i \|x_i - U U^T x_i\|^2 \quad \text{--- ①}$$

$$U^T U = I$$

Objective 2: Maximise the variance of the projected data

Mean-center the data, the resulting matrix is X . ($\mathbb{E}[X] = 0$)

$$\max_{U^T U = I} \mathbb{E} [\|U^T X\|^2]$$

or

For PCA, both objectives are equivalent

$$x = UU^T x + (I - UU^T)x$$

$$\|x\|^2 = \|UU^T x\|^2 + \|x - UU^T x\|^2 \quad (\text{cross-term} = 0)$$

$$\mathbb{E} [\|x\|^2] = \mathbb{E} [\|UU^T x\|^2] + \mathbb{E} [\|x - UU^T x\|^2]$$

$\left\{ U \text{ is a rotation, so } \mathbb{E} [U^T U] = \mathbb{E} [I] \right\}$

$$\mathbb{E} [\|x\|^2] = \mathbb{E} [\|U^T x\|^2] + \mathbb{E} [\|x - UU^T x\|^2]$$

fixed.

Hence, both objectives are equivalent ■

Finding U , consider objective 2.

Assume, we are projecting down to $k=1$.

$$\begin{aligned} \max_{\|u\|=1} & \frac{1}{n} \sum_i (u^T x_i)^2 \\ &= \max_{\|u\|=1} \frac{1}{n} \|U^T x\|^2 = \max_{\|u\|=1} u^T \underbrace{\left(\frac{1}{n} X X^T\right)}_{\text{covariance matrix (S)}} u = \text{largest eigenvalue of } S. \end{aligned}$$

Why?

$$\max_{\|u\|^2 \leq 1} u^T S u$$

$$L(u, \lambda) = u^T S u + \lambda (1 - u^T u)$$

$$\nabla_u L = 0 \Rightarrow 2S u - 2\lambda u = 0 \Rightarrow S u = \lambda u.$$

u is an eigenvector of S , with eigenvalue λ^1 .

$$\nabla_\lambda L = 0 \Rightarrow u^T u = 1.$$

$u^T S u = \lambda$, so pick u with largest eigenvalue

PCA:

Compute the covariance matrix S , and pick k eigenvectors of S with the highest eigenvalues

lecture

1. tensor puzzles (coding practice)

Bagging

Recall the bias variance decomposition for regression

$$\mathbb{E}[(y - h(x))^2] = \underbrace{\left(\bar{h}(x) - f(x)\right)^2}_{\text{Bias}^2} + \underbrace{\mathbb{E}[(h(x) - \bar{h}(x))^2]}_{\text{variance}} + \sigma^2$$

Consider constructing m independent datasets D_1, D_2, \dots, D_m all sampled from distribution P .

We can train a predictor for each of the n datasets $h(x; D_i)$

Average predictor $\frac{1}{m} \sum_{i=1}^m h(x; D_i)$

What happens to the bias and variance with using an average predictor.

$$\mathbb{E}_{D_1, \dots, D_m} \left[\frac{1}{m} \sum_{i=1}^m h(x; D_i) \right] = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{D_i, \dots, D_m} h(x; D_i) = \mathbb{E}_D h(x; D)$$

Variance

$$\text{Var}_{D_1, \dots, D_m} \left[\frac{1}{m} \sum_{i=1}^m h(x; D_i) \right] = \frac{1}{m^2} \sum_i \text{Var}_{D_i} [h(x; D_i)] = \frac{1}{m} \text{Var}_D (h(x; D))$$

Catch: We don't have m datasets, we only have one.

Solution: Bootstrap (Create multiple datasets from D using sampling with replacement)

Motivation for Bagging \rightarrow Bootstrap Aggregation

Bagging Algorithm

- Consider a dataset D with n training examples
- Create m different datasets, each of size n by sampling from D with replacement
- Average the predictions (for regression) and consider majority vote across predictions from models trained on each of the m datasets.

Note: Since the m datasets are not independent, we don't get the $\frac{1}{m}\sigma^2$ variance reduction. Can show that if the sampled predictions have a variance σ^2 , and a covariance of ρ^2 then,

$$\text{Var} \left[\frac{1}{m} \sum_i h(x; D_i) \right] = \frac{1}{m} (1-\rho) \sigma^2 + \rho \sigma^2$$

Each of these m datasets are referred to as Bootstrap samples.

A bootstrap sample contains roughly 63.2% of the observed training points in D .

Why? If we sample from a uniform distribution the prob. of a single dataset example not being drawn as a bootstrap sample datapoint is

$$P(\text{not drawn}) = \left(1 - \frac{1}{n}\right)^n$$

$$\text{As } n \rightarrow \infty = 1/e \approx 0.368.$$

The remaining instances in D that do not appear in a bootstrap sample are called "out of bag" samples [OOB]

OOB samples can be used to derive a bag estimate of test error

Process :

- For each i^{th} training example, predict the labels using all models trained using bootstrap sample for which i is OOB
- Aggregate the predicted values and compute the error for the i^{th} training eg.
- Average the error across all n training examples

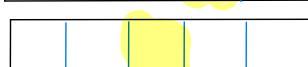
Random Forest

Decision trees + bagging + Tweak further to de-correlate trees

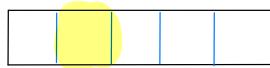
- Create bootstrap samples
- Train a decision tree for each bootstrap sample
- Tweak : For every split in every DT, only consider a subset of features.
If there are d features, we lose a random subset of size k (where typically, $k \approx \sqrt{d}$) at each split for every DT.

Cross-validation

Technique used for model evaluation / selection.



D



test train

(k-fold) ← partition

$k=5$

leave-one-out
cross validation
($k=n$)

Typically $k=3, 5, 10$.

Note: Every example shows up in both training, test at some point