

Machine Learning

Supervised Learning:

$$\text{Training set} \rightarrow S_n = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

$$x_i \in X$$

$$y_i \in Y$$

$$\hat{f}: X \rightarrow Y$$

$$S_r \subseteq \underbrace{X \times Y}_{\text{data space}}$$

y_i ↗ regression
classification
 $\{-1, 1\}$

$\hat{f}(x_i) \approx y_i$ ability to represent the available data, we want our function to generalize to new unseen data.

Goal: given a training set, Learn input-output relation

assumptions $p(x, y) = p_x(x)p(y|x)$

Joint-probability distribution
Unknown

marginal distribution
uncertainty

foreach x , there is possible output $p(y|x)$

in Regression

$$y = f(x) + \varepsilon \quad \begin{matrix} \text{random noise eq gaussian} \\ \text{regression function} \end{matrix}$$

We assume \hat{f} will generalize well to future data

in Classification

$$p(1|x) = 1 - p(-1|x)$$

LOSS FUNCTION quantifies how good is the prediction wrt real data

$l(y, f(x))$ point-wise error measure

* how much we lose if we use the prediction with estimation $\hat{f}(x_i)$

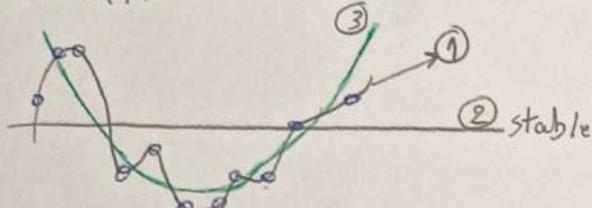
Expected Loss (expected risk)

joint probability not y_i

$$\mathcal{E}(f) = \mathbb{E}[l(y, f(x))] = \int p(x, y) \cdot l(y, f(x)) dx dy$$

best function = target function = f^* that minimizes expected risk.

$$f^* = \operatorname{argmin} E(f) \text{ one that minimize } E$$



- function:
- ① error=0 → too complex
→ not stable → overfitting
 - ② too stable → underfitting
 - ③ compromise

\hat{f}_s : depend on Training set

excess expected loss

↓ we want it small

$$E(f^*) - E(f_s)$$

$$\mathbb{E}_s [|E(f^*) - E(f_s)|]$$

$$\text{consistency } \lim \mathbb{E}_s [|E(f^*) - E(\hat{f}_s)|] = 0$$

Good algorithm: ① Fitting: Fit data well ② Stability: should not change much if data change slightly

Regularization Parameter: controls tradeoff between data-fitting & stability

Empirical Risk

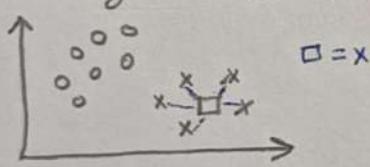
$$\hat{E}(f) = \frac{1}{n} \sum_{i=1}^n l(y_i, f(x_i))$$

Empirical Risk minimization $\Rightarrow \hat{f} = \arg \min \hat{E}(f)$

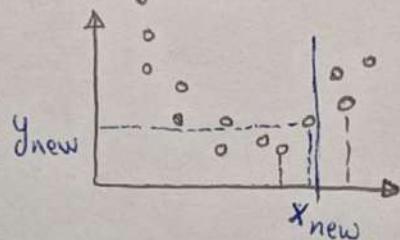
Local Methods

* close or near points in input $\xrightarrow{\text{expected}}$ close or near in output space

Binary classifi.



Regression



we assign to output of x_{new} same y value of closest point

NN: Nearest neighbor

$$y_{\text{new}} = \text{Nearest Neighbor}(x_{\text{new}}, X_n, Y_n)$$

$$X_n \times d$$

$$d = \text{zeros}(1, n)$$

$$Y_n \times 1$$

$$\text{for } i = 0 \dots n-1:$$

$$d(i) = \|x_{\text{new}} - x_n(i, :)\|$$

$$i_{\text{min}} = \operatorname{argmin}(d)$$

$$Y_{\text{new}} = Y_n(i_{\text{min}})$$

K-NN

Binary Classification

$$Y_{\text{new}} = \text{KNN}(x_{\text{new}}, X_m, Y_m, K)$$

$$d = \text{zeros}(1, m)$$

$$O(nd) \quad \begin{bmatrix} \text{for } i=0 \dots m-1 \\ d(i) = \|x_{\text{new}} - X_m(i, :) \| \end{bmatrix}$$

$$O(n \log n) \quad [d_{\text{sort}}, i_{\text{sort}}] = \text{argsort}(d, \text{ascend})$$

$$O(1) \quad Y_{\text{new}} = \text{sign} \left(\sum (Y_m(i_{\text{sort}}^{0 \rightarrow K-1})) \right)$$

Regression

$$Y_{\text{new}} = \text{KNN}(X_{\text{new}}, X_m, Y_m, K)$$

$$d = \text{zeros}(1, m)$$

for $i=0 \dots m-1$

$$d(i) = \|x_{\text{new}} - X_m(i, :) \|$$

$$[d_{\text{sort}}, i_{\text{sort}}] = \text{argsort}(d, \text{asc})$$

$$Y_{\text{new}} = \text{average}(Y_m(i_{\text{sort}}^{0 \rightarrow K-1}))$$

KNN

$$f(x) = \frac{1}{K} \sum_{i=1}^K y_{j_i}$$

$K = \text{odd}$

bigger $K \rightarrow$ more stable

KNN depends on K

small $K \rightarrow$ overfitting

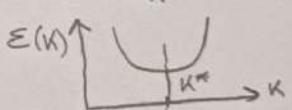
very large $K \rightarrow$ very stable, disregard data
underfitting

$f_{s, K}$ → hyperparameter of the method
behaviour of alg depend on K

best/optimal choice for hyperparameter K : The expected Loss \mathcal{E}_K

$$\mathcal{E}_K = \mathbb{E}_s \mathbb{E}_{x, y} (y - \hat{f}_{s, K}(x))^2$$

$$K^* = \operatorname{argmin} \mathcal{E}_K$$



K^* minimize \mathcal{E}_K

NB: There is always some noise

$$\mathcal{E}(K) = \mathbb{E}_s \left(\underbrace{\hat{f}(x) - \hat{f}_{sK}(x)}_{\text{error term}} \right)^2 + \underbrace{\sigma^2}_{\text{noise}}$$

Idealized K-NN Function

$$\hat{f}_{sK}(x) = \frac{1}{K} \sum f^*(x_l)$$

$$\text{Standard KNN func } \hat{f}_{sK}(x) = \frac{1}{K} \sum y_l = f^*(x_l) + \underbrace{\delta}_{\text{noise}}$$

adding & subtracting to $\mathbb{E}(K)$ & simplifying ...

$$\mathbb{E}(K) = \underbrace{\left[f^*(x) - \mathbb{E}_{\tilde{f}_{SK}}[\tilde{f}(x)] \right]^2}_{\text{bias}} + \underbrace{\mathbb{E}_s \left[\tilde{f}_{SK}(x) - \hat{f}_{SK}(x) \right]^2 + \sigma^2}_{\text{variance}}$$

$\mathbb{E}(K) = \text{bias} + \text{variance}$

variance: $\mathbb{E}_s \left(\tilde{f}_{SK}(x) - \hat{f}_{SK}(x) \right)^2 + \sigma^2$

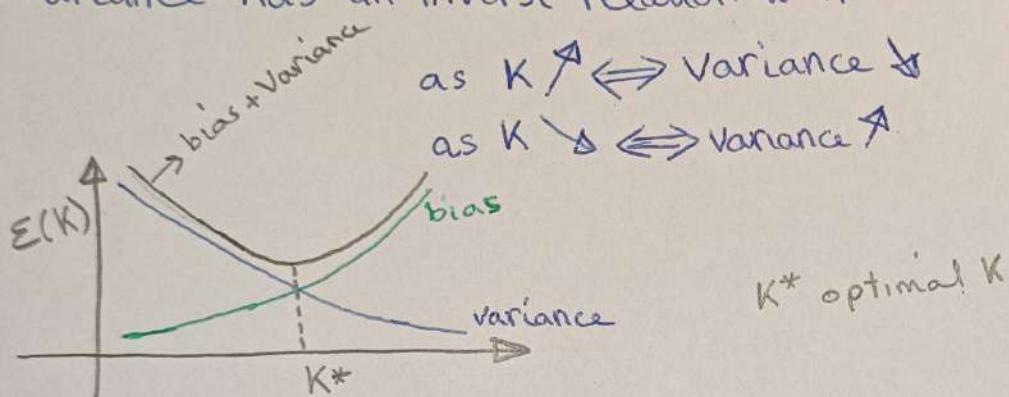
$$= \mathbb{E}_s \left[\frac{1}{K} \sum_{l \in Kx} f^*(x_l) - \frac{1}{K} \sum_{l \in Kx} y_l \right]^2 + \sigma^2 = \frac{1}{K^2} \mathbb{E}_s \left[\sum_{l \in Kx} (f^*(x_l) - y_l)^2 \right] + \sigma^2$$
$$= \dots = \frac{\sigma^2}{K} + \sigma^2$$

* Variance has an inverse relation w/ K

as $K \nearrow \Leftrightarrow$ Variance \downarrow

as $K \searrow \Leftrightarrow$ Variance \uparrow

variance depend on
K and noise



$E(K)$ compromise
between bias & variance

Notes : ChatGPT

Underfitting \rightarrow high bias

Overfitting \rightarrow high Variance

Bias : * approximating a real-world problem which may be complex by a simple model

* high bias \rightarrow underfitting \rightarrow unable to capture patterns of data

Variance : * represents model sensitivity to small fluctuations/noise

* high variance \rightarrow overfitting \rightarrow model too complex and fitting training data well, but capturing noise

* overfitting \rightarrow perform well on training data but fail to generalize to new unseen data

Goal: right balance

Bias \rightarrow fitting

Variance \rightarrow stability

Cross Validation train on some data and validate on new unseen data

To choose K in practice

T	V
---	---

✓ usually randomly sampled from dataset

$K = 1, 3, \dots, K_{\max} \rightarrow$ foreach K , execute procedure and estimate error

* $K^* = \text{cross-validation}(T, V, K_{\max})$

↓
optimal
 K

$i = 0$
for $K = 1, 2, \dots, K_{\max}$

- get $\hat{f}_{T,K}$
- error[i] \leftarrow error $\hat{f}_{T,K}$ on V
- $i++$

return K // such that error is minimum

V	T
---	---

T	V	T
---	---	---

T	V	T
---	---	---

,

* For each K

① Shuffle and split the data into T (training) & V (validation)

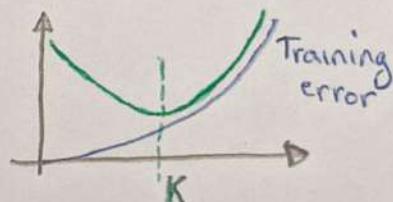
② Train algorithm on T

③ compute empirical loss on V

$$\hat{\epsilon}_K = \frac{1}{|V|} \sum_{x,y \in V} (y - \hat{f}_{T,K}(x))^2$$

we work on V

④ Select \hat{K} that minimize $\hat{\epsilon}_K$ because we compute on Validation set



* $K\text{-Fold} = 5$

$$\text{error} = \frac{1}{5} \sum_{i=1}^5 \text{err}[i]$$

$V_1 \rightarrow \text{err}[1]$

$V_2 \rightarrow \text{err}[2]$

⋮
 $V_5 \rightarrow \text{err}[5]$

Linear Models

ERM: empirical error $\rightarrow \min_{w \in \mathbb{R}^d} \hat{L}(w)$

General Idea: $\hat{\Sigma}(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$

$$= \min \sum_{i=1}^n (y_i - \omega^T x_i)^2$$

proxy for expected error

$$\Sigma(f) = \mathbb{E}[\ell(y, f(x))] = \int p(x, y) \ell(y, f(x)) dx dy$$

is unknown

* $\omega^T x + b = \sum_{i=1}^n \omega_i x_i + b$

Empirical Risk:
 $\min \hat{L}(w)$

function f
defined by
vector w

$$\hat{L}(w) = \sum_{i=1}^n (\omega^T x_i - y_i)^2$$

(NB)

* NN is local

* Linear is global \rightarrow new value
even if far can change model

$$\nabla \hat{L}(w) = \sum_{i=1}^n 2(\omega^T x_i - y_i) x_i = 0$$

x_i is row of \hat{X} matrix

* $L(w) = \|\hat{X}\omega - \hat{Y}\|^2$

$\nabla \rightarrow$ gradient \rightarrow allows to find w
direction & rate of
steepest increase in
Loss function

$$\Rightarrow \nabla \hat{L}(w) = 2 \hat{X}^T (\hat{X}\omega - \hat{Y}) = 0$$

$$\hat{X}^T \hat{X} \omega = \hat{X}^T \hat{Y}$$

$$\omega = \frac{\hat{X}^T \hat{Y}}{\hat{X}^T \hat{X}}$$

must be
invertible

$$\hat{X}^T \hat{X} = V S V^T \quad v \begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix} v^T \begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix} = I$$

$n > d$: Underparameterized (over determined)

building of $\hat{X}^T \hat{X}$ complexity

$$\mathcal{O}(d^2 n + d^3)$$

multiplic. inverting

$$\hat{X}^T \omega / \hat{X}$$

$$\frac{\hat{X}^T \hat{X}}{d \times d} \omega = \frac{\hat{X}^T \hat{Y}}{d \times 1}$$

$$\mathcal{O}(d^2 n + d^3)$$

$d > n$: Overparameterized building $XX^T \rightarrow \mathcal{O}(n^2 d + n^3)$

$$\hat{X} \hat{X}^T c = \hat{Y}$$

$$c = \frac{\hat{Y}}{\hat{X} \hat{X}^T}$$

$$\omega = \hat{X}^T c$$

building $XX^T \rightarrow \mathcal{O}(n^2 d + d^3)$

$$\omega = \hat{X}^T (\hat{X} \hat{X}^T)^{-1} \hat{Y}$$

Pseudo Code

$$\omega = LS_Train(\hat{X}, \hat{Y})$$

$$n = X.shape[0]$$

$$d = X.shape[1]$$

if $n \geq d$:

$$\omega = \hat{X}^T \hat{Y} / \hat{X}^T \hat{X}$$

else:

$$c = \hat{Y} / \hat{X} \hat{X}^T$$

$$\omega = \hat{X}^T c$$

TIKHONOV REGULARIZATION

$\lambda > 0$ very small

$$\hat{X}\hat{X}^T = U \begin{pmatrix} \frac{1}{\sigma_1} & & \\ & \ddots & \\ & & \frac{1}{\sigma_n} \end{pmatrix} U^T$$

- * if λ is too big: matrix becomes (0) and
 - ① $y_i = \omega^T x = 0$
 - ② σ_i insignificant
 - all eigenvalues $\frac{1}{\lambda}$
 - matrix scales input

add small amount to denominator

$$\left(\frac{1}{\sigma_i + \lambda} - \frac{1}{\sigma_i} \right)$$

$$\begin{array}{ll} \omega \xrightarrow{d > n} \lambda = 0 & \omega = \hat{X}^T (\hat{X}\hat{X}^T)^{-1} \hat{Y} \\ & \lambda \neq 0 & \omega = \hat{X}^T (\hat{X}\hat{X}^T + \lambda I)^{-1} \hat{Y} \\ \xrightarrow{n > d} \lambda = 0 & \omega = (\hat{X}^T \hat{X})^{-1} \hat{X}^T \hat{Y} \\ & \lambda \neq 0 & \omega = (\hat{X}^T \hat{X} + \lambda I)^{-1} \hat{X}^T \hat{Y} \end{array}$$

Regularized ERM

$$* L_\lambda(\omega) = \sum_{i=1}^n (y_i - \omega^T x_i)^2 + \underbrace{\lambda \|\omega\|^2}_{\text{penalty regularizer}} \Leftrightarrow \|\hat{X}\omega - \hat{Y}\|^2 + \lambda \|\omega\|^2$$

adding $\lambda \|\omega\|^2$ is called RIDGE REGRESSION

$$\nabla L_\lambda(\omega) = 2\hat{X}^T (\hat{X}\omega - \hat{Y}) + 2\lambda \omega = 0$$

$$(\hat{X}^T \hat{X} + \lambda I) \omega = \hat{X}^T \hat{Y}$$

Empirical Risk Minimization: $\min \hat{L}(\omega)$

Regularizer controls \rightarrow stability
 \rightarrow prevent overfitting

$\|\omega\|^2$: regularizer

λ balances error term & regularizer

RLS: Regularized Least Square $\min_{\omega^*} \hat{L}(\omega)$

$$\min_{\omega \in \mathbb{R}^D} \frac{1}{n} \sum_{i=1}^n (y_i - \omega^T x_i)^2 + \lambda \|\omega\|^2$$

LOGISTIC REGRESSION

* Logistic Loss Function $\ell(y, f_w(x)) = \log(1 + e^{-y f_w(x)})$

y not binary

$$\frac{d}{da} \log(1 + e^{-a}) p(1|x) + \log(1 + e^a) p(-1|x)$$

$$= \frac{-e^{-a}}{1 + e^{-a}} p(1|x) + \frac{e^a}{1 + e^a} p(-1|x)$$

$$= \frac{-1}{1 + e^a} p(1|x) + \frac{e^a}{1 + e^a} p(-1|x) = 0$$

$$-p(1|x) + e^a p(-1|x) = 0$$

$$e^a = \frac{p(1|x)}{p(-1|x)} \neq$$

$$\Leftrightarrow a_* = \log \frac{p(1|x)}{p(-1|x)}$$

$$\begin{aligned} \log(>1) &\Leftrightarrow a > 0 \\ \log(<1) &\Leftrightarrow a < 0 \end{aligned}$$

log is monotonic increasing

$$\therefore a > 0 \rightarrow p(1|x) > p(-1|x)$$

$$a < 0 \rightarrow p(1|x) < p(-1|x)$$

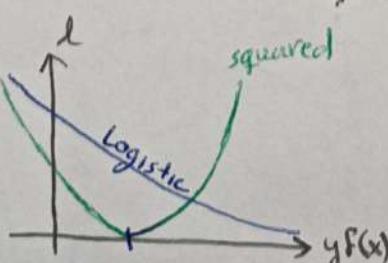
$$\neq e^{a_*} = \frac{p(1|x)}{p(-1|x)}$$

$$e^{a_*} = \frac{p(1|x)}{1 - p(1|x)}$$

$$p(1|x) = e^{a_*} (1 - p(1|x))$$

$$= \frac{1}{1 + e^{-a_*}}$$

$$p(1|x) + p(-1|x) = 1$$



Binary $y \in \{-1, 1\}$

$$L(f) = \min_{f(x)} L(f) = \int l(f(x), y) p(x, y) dx dy$$

$$= \int (\ell(f(x), y) p(y|x) dy) p(x) dx$$

$$\text{Let } L(f(x)) = L(a)$$

$$L(a) = \int \ell(a, y) p(y|x) dy$$

$$= \ell(a, 1) p(1|x) + \ell(a, -1) p(-1|x)$$

derive $L(a)$ where $l = (y - a)^2$

$$\frac{d}{da} [(1-a)^2 p(1|x) + (-1-a)^2 p(-1|x)] = 0$$

$$\Rightarrow -2() p(1|x) - 2() p(-1|x) = 0$$

$$\Rightarrow -a(p(1|x) + p(-1|x)) + p(1|x) - p(-1|x) = 0$$

$$\neq \Rightarrow a_* = p(1|x) - p(-1|x)$$

optimal $a = a_*$

$$f(x) = a_* = \begin{cases} > 0 & p(1|x) > p(-1|x) \\ < 0 & p(-1|x) > p(1|x) \end{cases}$$

if $p(1|x) > p(-1|x)$ makes sense to predict 1

$$\neq a_* = p(1|x) - p(-1|x)$$

add $p(1|x)$ on both sides

$$a_* + p(1|x) = 2p(1|x) - p(-1|x)$$

$$\begin{aligned} p(1|x) &= \frac{a_* + p(1|x) + p(-1|x)}{2} \\ &= \frac{a_* + 1}{2} \end{aligned}$$

Logistic Regression Minimizer

Gradient Descent

$$L_{\lambda}^{(w)} = \frac{1}{n} \sum_{i=1}^n \log \left(1 + e^{-y_i w^T x_i} \right) + \frac{\lambda \|w\|^2}{2}$$

$$\nabla L_{\lambda}^{(w)} = \frac{1}{n} \sum_{i=1}^n \underbrace{\frac{e^{-y_i w^T x_i} \cdot (-y_i x_i)}{1 + e^{-y_i w^T x_i}}}_{\alpha} + 2\lambda w = 0$$

* $w_{t+1} = w_t - \gamma F'(w_t)$

ex $w_1 = w_0 - \gamma F'(w_0)$

$\alpha = F'(w)$

$t = \text{iteration}$

$$f = \log e^u = \frac{u}{u}$$

$$f' = \frac{u'}{u}$$

to memorize
 $w^T w \rightarrow 2w$

$$-y_i x_i w^T = aw^T$$

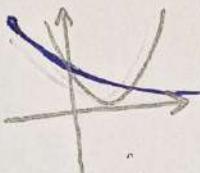
* use Newton's law to know if we are close to minimum??

$$w_{t+1} = w_t - \frac{F'(w_t)}{F''(w_t)}$$

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

$$\text{and } \gamma_t = \frac{1}{F''(w_t)} = \left(F''(w_t) \right)^{-1}$$

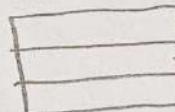
* GD goes down by e^{-t}



squared is simpler than logistic

* Gradient descent: * is an optimization alg. * goal to minimize Loss function * goal to find parameters (weight-)

optimization



compute
 ∇L in pieces

Newton's method is more computational costly

* gradient provides direction of steepest Ascent

* γ : Learning space \rightarrow positive scalar \rightarrow controls the size of the steps taken

Stochastic Gradient descent SGD used in Large Datasets

- * instead of computing gradient for entire dataset
 - \rightarrow compute gradient of loss for single point or small random subset
- * frequent update \rightarrow faster convergence
- * model parameters are updated after processing each data point or mini-batch

$L(w_t; x_i, y_i)$: Loss
 function evaluated at i -th data point for w

$$w_{t+1} = w_t + \gamma_t \nabla L(w_t; x_i, y_i)$$

Variations of gradient descent

- mini-batch
- Stochastic GD
- SGD

- Adam

SGD slower
 goes down by $\frac{1}{\sqrt{t}}$

Non-Linear Models

$$\hat{z} = \text{feature_map}(\hat{x}, \text{"type"}, \text{parameter})$$

Linear Model: $\omega^T x_i = y_i$ $i = 1 \dots n$

Non-Linear: such as $w^T x^2$

• $f(x) = ax^2 + bx = w^T z$

$$z = (x^2, x) \text{ and } \omega = (a, b)$$

$$\hat{z} = \text{feature_map}(\hat{x}, \text{"poly", "degree": 2})$$

$$i = 1 \dots n$$

$$z_i = x_i^2$$

• $f(x) = ax^3 + bx^2 + cx + d$
 $= w^T z$

$$\omega = (a, b, c, d)$$

$$z = (x^3, x^2, x, 1)$$

• superposition of sin & cos $\sum_{j=1}^P c_j \sin(2\pi j x) = f(x) = w^T z$
 nb of frequencies

$$\omega^j = c^j \quad j = 1 \dots P$$

$$z^j = \sin(2\pi j x)$$

• gaussian $\sum_{j=1}^P e^{-\frac{|x - m_j|^2}{2\sigma^2}} = f(x)$
 $m_j = \text{mean}$
 $\sigma = \text{std. dev}$

all models
can be written
as:

$$f(x) = \sum_{j=1}^P c_j \varphi_j(x)$$

φ : small phi

$$\varphi_j: x \rightarrow \mathbb{R} \quad j = 1 \dots P$$

$$x_i \rightsquigarrow z_i = (\varphi_1(x_i), \dots, \varphi_P(x_i))$$

feature

$$z = (\varphi_1(x) \dots \varphi_P(x)) = \Phi(x)$$

and $f(x) = w^T \Phi(x) = \sum_{j=1}^P w^j \varphi_j(x)$

• $w = (\hat{\Phi}^T \hat{\Phi} + n\lambda I)^{-1} \hat{\Phi}^T \hat{Y}$ $O(P^2 n + P^3)$

$$w = \Phi^T \underbrace{(\hat{\Phi} \hat{\Phi}^T + n\lambda I)^{-1}}_C \hat{Y} \quad O(n^2 P + n^3)$$

$$w = \Phi^T c$$

$$w = (\hat{x}^T \hat{x} + n\lambda I)^{-1} \hat{x}^T \hat{y}$$

$$w = \hat{x}^T (\hat{x} \hat{x}^T + n\lambda I)^{-1} \hat{y}$$

Scales exponentially \rightarrow too many coefficients
 too many params
 \rightarrow over parameterization

We choose this bcz P very big

Non-Linear Models

$$\hat{z} = \text{feature_map}(\hat{x}, \text{"type"}, \text{parameter})$$

Linear Model: $\omega^T x_i = y_i$ $i = 1 \dots n$

Non-Linear: such as $w^T x^2$

* $f(x) = ax^2 + bx = w^T z$

$$z = (x^2, x) \text{ and } w = (a, b)$$

* $\hat{z} = \text{feature_map}(\hat{x}, \text{"poly", degree=2})$

$$i = 1 \dots n$$

$$z_i = x_i^2$$

* $f(x) = ax^3 + bx^2 + cx + d$
 $= w^T z$

$$w = (a, b, c, d)$$

$$z = (x^3, x^2, x, 1)$$

* $\hat{z} = \text{feature_map}(\hat{x}, \text{"poly", degree=3})$

$$i = 1 \dots n$$

$$z_i = (x_i^3, x_i^2, x_i, 1)$$

* superposition of sin & cos $\sum_{j=1}^P c_j \sin(2\pi j x) = f(x) = w^T z$
 nb of frequencies

$$w^j = c^j \quad j = 1 \dots P$$

$$z^j = \sin(2\pi j x)$$

* gaussian $\sum_{j=1}^P e^{-\frac{|x - m_j|^2}{2\sigma^2}} = f(x)$

m_j = mean

σ = std. dev

all models
can be written
as:

$$f(x) = \sum_{j=1}^P c_j \varphi_j(x)$$

φ : small phi

$$\varphi_j: x \rightarrow \mathbb{R} \quad j = 1 \dots P$$

$$x_i \rightsquigarrow z_i = (\varphi_1(x_i), \dots, \varphi_P(x_i))$$

feature

$$z = (\varphi_1(x) \dots \varphi_P(x)) = \Phi(x)$$

and $f(x) = w^T \Phi(x) = \sum_{j=1}^P w^j \varphi_j(x)$

* $w = (\hat{\Phi}^T \hat{\Phi} + n\lambda I)^{-1} \hat{\Phi}^T \hat{Y}$ $O(P^2 n + P^3)$

$$w = \Phi^T \underbrace{(\hat{\Phi}^T \hat{\Phi} + n\lambda I)}_C \hat{Y} \quad O(n^2 p + n^3)$$

$$w = \Phi^T c$$

$$w = (\hat{x}^T \hat{x} + n\lambda I)^{-1} \hat{x}^T \hat{Y}$$

$$w = \hat{x}^T (\hat{x} \hat{x}^T + n\lambda I)^{-1} \hat{Y}$$

Scales exponentially \rightarrow too many coefficients

+ too many params

\rightarrow overparameterization

We choose this bc P very big

Kernels

Representer Theorem:

$$f(x) = \sum_{i=1}^n \Phi^T(x_i) \Phi(x) c_i$$

$$K(x_i, x) = \Phi^T(x_i) \Phi(x)$$

KRLS_PREDICT

$$f(x) = \sum_{i=1}^n K(x_i, x) c_i$$

also $\omega = \Phi^T(\Phi\Phi^T + n\lambda I)\hat{Y} \Rightarrow \omega = \Phi^T c$

$$c = (\Phi\Phi^T + n\lambda I)\hat{Y}$$

KRLS_TRAIN

$$c = (\hat{K} + n\lambda I)\hat{Y}$$

Kernels

① Linear:

$$K(x, x') = x^T x'$$

(ex) or

$$K(x_i, x_j) = x_i^T \cdot x_j$$

$$K(x, z) = x^T \cdot z$$

② Polynomial: $K(x, x') = (x^T x' + 1)^d$

d = degree

(ex) or

$$K(x, z) = (xz + 1)^p$$

$$K(x_i, x_j) = (x_i^T x_j + 1)^d$$

③ Gaussian:

$$K(x, x') = e^{-\frac{\|x - x'\|^2}{2\sigma^2}}$$

(ex)

$$K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

Kernel: ① Symmetric $\rightarrow K(x, z) = K(z, x)$ for all $x, z \in \mathbb{R}^d$

② positive definite (p.d.)

Neural Networks

Recap Linear $f_w(x) = w^T x \quad x \in \mathbb{R}^d$

best vector $w^* = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n l(y_i, w^T x_i) + \lambda \|w\|^2$

for Logistic regression, we use iterative method based on gradient descent

$$w_0, w_t = w_{t-1} + \gamma \nabla G(w)$$

\downarrow step

non-Linear $\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^P$ feature map to transform data

to linear, adding new dimensions that allow us to create hyperplane that separate data

$$w^* = \underset{w \in \mathbb{R}^P}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n l(y_i, w^T \Phi(x_i)) + \lambda \|w\|^2 \quad \text{solution now lives in } P \text{ dimension}$$

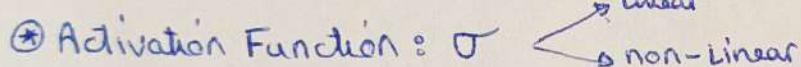
* P very large \rightarrow Kernel trick offer efficient approach

$$f(x) = \frac{1}{n} \sum_{i=1}^n c_i K(x, x_i)$$

solution is now collection of coefficients
 $c_i \rightarrow$ we have n coeffs.

Neural Networks

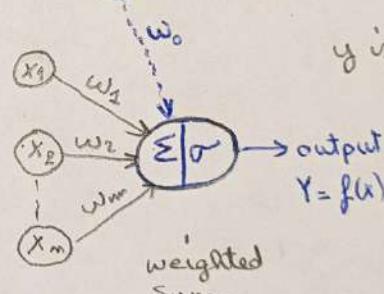
* linearly combines input

* Activation Function: σ 

* Single Layer: $y = f(x)$

$$f(x) = \sigma \left(\sum_{i=1}^m w_i x_i \right)$$

$$\text{if } w_0 \rightarrow \text{or } f(x) = \sigma \left(\sum_i w_i x_i + w_0 \right)$$



y is a binary signal

$$g(x) = \sum_{i=1}^m x_i \quad \theta = \text{threshold}$$

$$f(x) = \sigma(g(x)) = \begin{cases} 1 & g(x) > \theta \\ 0 & \text{otherwise} \end{cases}$$

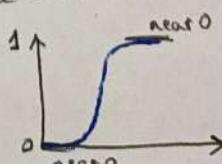
σ : activation func

* Activation Functions: - allow to account for non-linearity in model
 - take single nb and map it to different numerical value

(1) Sigmoid: - corresponds to Logistic Regression

- non-zero centered output
- suffer from vanishing gradient
 - * issue of vanishing gradient when the gradient becomes very small during back-propagation & weights of network updated very little and learning becomes slow

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

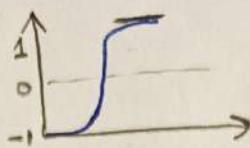


recap
 → gradient of loss function is a vector that points in the direction of steepest increase of loss
 → gradient provide info of how loss would change if we make small adjustments to param

- ② Tanh:
- * it is a scaled version of sigmoid
 - * suffers from vanishing gradient problem
 - * output is zero centered (better gradient prop.)

The tanh function maps input x to output in range $[-1, 1]$

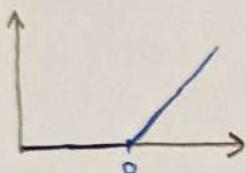
$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



scaled sigmoid = $2 \sigma(2x) - 1$

- ③ ReLU (Rectified Linear Unit) good for images

$$f(x) = \max(0, x)$$



- * ReLU activation funct. returns x for any positive input and zero for any negative input. Looks linear for positive & flat for zero.
- * it is simple to compute, making it computationally efficient during both forward & backward propagation
- * ReLU does not saturate for positive values \rightarrow which helps in vanishing gradient problem \rightarrow accelerate training
- * weight may irreversibly die

- ④ Dying ReLU: relu becomes inactive (zero output) for all inputs during training and fail to update weights \rightarrow issue occur due to large gradient flow. ReLU network \rightarrow cause to output zero in subsequent updates
Solution: Leaky ReLU

- ④ Leaky ReLU:

$$f(x) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha x & \text{if } x < 0 \end{cases}$$

- * fix dying ReLU problem
- * introduce small, non-zero slope for negative inputs, allowing the pass of info when input is negative
- * prevent neurons from being inactive by allowing small, non-zero gradient for negative input during backward pass

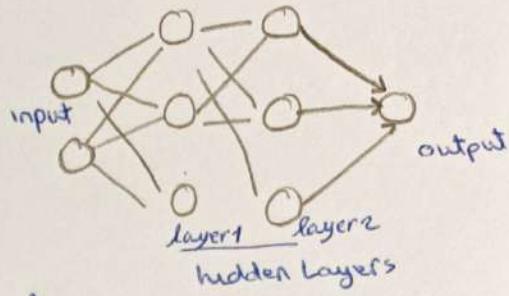
α is a small positive constant that determine slope of function for negative input

Activation functions so important because they introduce non-linearity into the network.

Multilayer networks

chain of Layers

- Input Layer
- Multiple hidden Layers
- Output Layer



- * each element of Layer \rightarrow neuron ; width = dimensionality of hidden layers
- ↙
each neuron have activation function

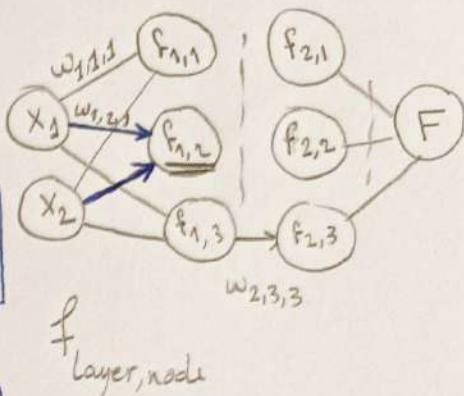
take $f_{1,2}$

$$f_{1,2} = \sigma \left[\left(\sum_{i=1}^2 w_{2,2,i} \cdot x_i \right) + b_{1,2} \right]$$

take $f_{2,1}$

$$f_{2,1} = \sigma \left[\left(\sum_{i=1}^3 w_{2,1,i} \cdot f_{1,i} \right) + b_{2,1} \right]$$

$F(x) = \sigma \left(\underbrace{w_{3,1} \cdot f_{2,1}} + \underbrace{w_{3,1,2} \cdot f_{2,2}} + \underbrace{w_{3,1,3} \cdot f_{2,3}} + b_{3,1} \right)$



f
layer, node

w
layer, node, input

④ Optimizing a neural network: determining the weights that are most appropriate for our problem solution

④ different layers can have different Activation function

④ same Layer should have same Activation function

Training DNN

means learning the values for the model parameters (w , bias)
from Training Set \rightarrow use loss func
to estimate how much we loose predicting $F(x)$

How many parameters

○ ○
○ ○ ○ $\sigma F(x)$
○ ○ ○

18 weights

+ 7 biases

= 25

input
no bias
unknowns

Problem: the function we want to minimize is very complicated \rightarrow difficult for gradient descent.

computing gradient

$$w^* = \operatorname{argmin} \frac{1}{n} \sum_{i=1}^N l(F(x_i; w), y_i)$$

for Complex! function

$$w^* = \operatorname{argmin} J(w)$$

use gradient descent: $w_{t+1} = w_t + \gamma \nabla J(w_t)$

gradient = 0
 $\Rightarrow J(w)$ minimal
 w^* when gradient = 0

Forward propagation: input flow into network & produce cost

⊗ Backward propagation: info flows back the net to compute gradient

Backpropagation
minimize cost func by adjusting weights & biases
level of adjustment determined by gradients of cost func.
goal is to compute partial derivatives wrt any w or b.

entire dataset too big → we use mini-batch to train

$$\omega_{t+1} = \omega_t + \gamma \nabla J(\omega)$$

Reverse Model

$$\frac{dF}{dx} = \frac{dF(x)}{dx} = \frac{df}{dx} f(g(h(u(v(x)))))$$

chain rule is essential part of DNN:

↳ goal: estimate the gradient

for high dimensional spaces wrt each weight

start from outer function (right to left)

Gradient Descent Recap

- mini-batch

- stochastic

more computationally efficient
more iterations
smaller steps
more unstable because in some iterations might go opposite direction

DNN: dense NN

⊗ How does a small change in weight affect loss $J(W)$

↳ we compute gradient with respect to the nb. of w's

Neural Networks → have lot of Local minima → very difficult to find global minimum

* test, validation error against nb epochs

* we can add regularization term (λ) → cross-validation too costly - not feasible

Summary: we want to find parameters that minimize cost function

$$w^* = \underset{w}{\operatorname{argmin}} \frac{1}{n} l(\hat{y}_i, y_i), \quad y_i = F(x_i, w)$$

Notes

DNNs → high nb of parameters → high variance → high chance of overfitting
high Variance → overfitting → poor response to noise → does not generalize well
other notes ⊗ high bias → underfitting

Terms: epoch: when entire dataset sent forward & backward through neural network only once

batch: training examples in a mini-batch

Recap: in DNN, we have lot of parameters → can result in overfitting
→ not well generalize
→ poor response to noise

possible solution: Regularization to stabilize learning weights

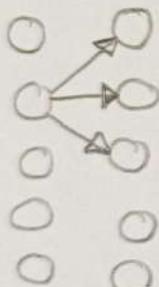
- balanced ✓
 - strong regularization → add significant bias → underfitting
 - weak regularization → not able to capture patterns
- Lead to high variance → overfitting

Convolutional Neural Networks

CNN: special kind of NN with a known grid-like topology

CNN

2 properties: ① Sparse interaction



NN: output of node is input to all next layer nodes

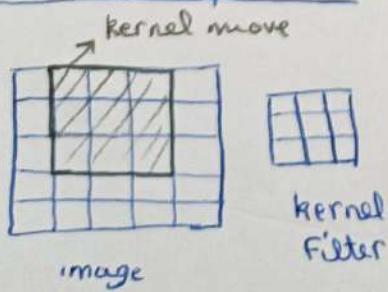
CNN: output of node ^{is} input to some neurons in next layer or input comes from a subset of nodes

positives of few parameters → Lower memory requirement
computations → statistical efficiency
reduced variance when estimating parameters

② Parameter sharing:

in CNN: weight used at every position of the input

Convolution Operation



(m, n) : indices used in summation
↳ dimensions of kernel

Convolution between image & kernel

$$s(i, j) = (K * I)(i, j) = \sum_m \sum_n I(i+m, j+n) K(m, n)$$

$s(i, j)$ → value of output (feature map) at position (i, j)

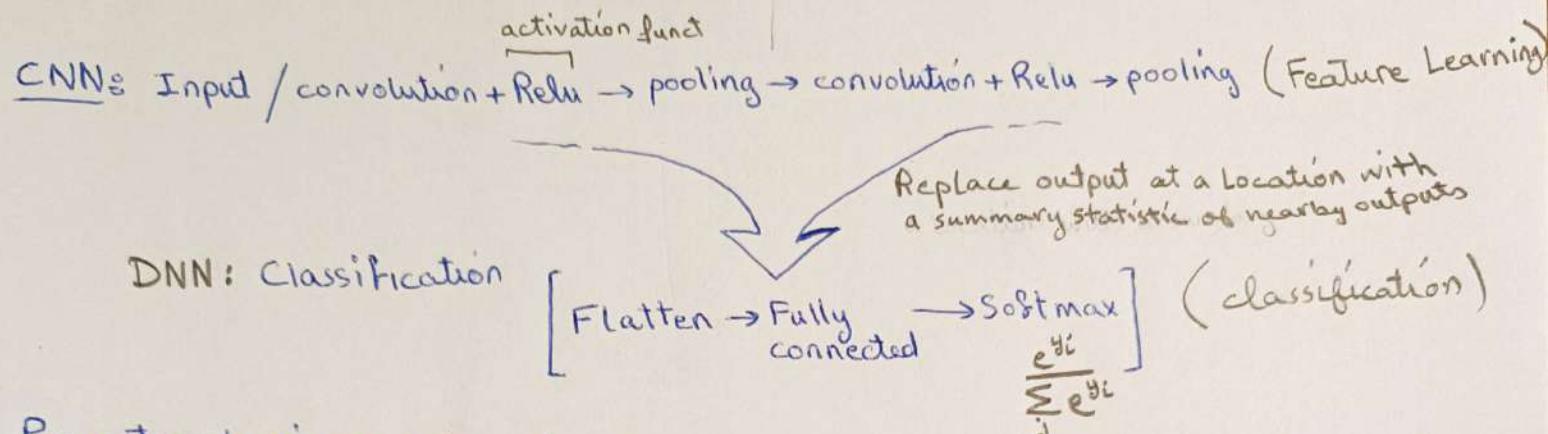
$(K * I)(i, j)$ → convolution operation b/w Kernel / Filter K and the input at position (i, j)

$K(m, n)$ → kernel value at position (m, n)

$I(i+m, j+n)$ → Image value at position $(i+m, j+n)$

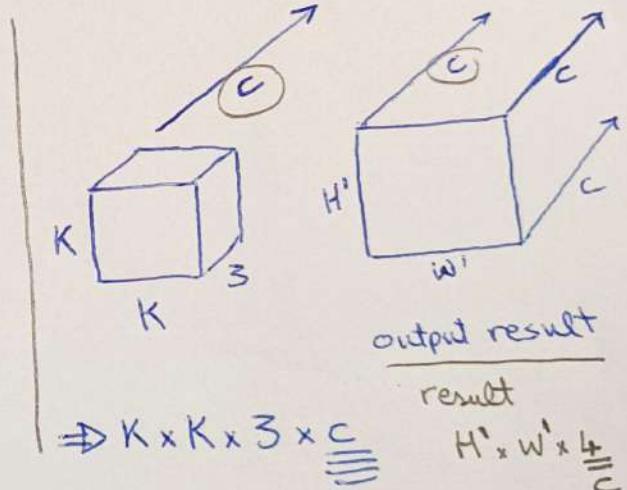
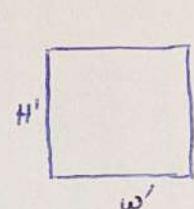
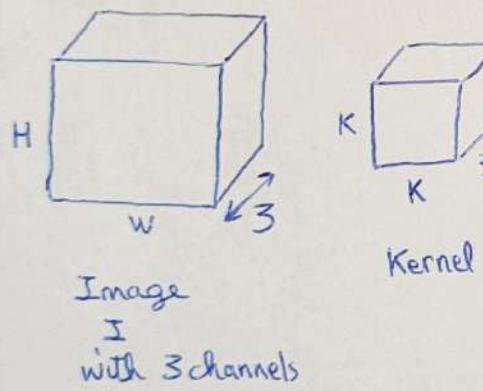
- CNN:
- * nowadays, learn features using CNN
 - * Kernel: is a feature detector slides over inputs to generate Feature Map.

[ex] image $200 \times 200 \times 3 = \text{size} = 120000$ weights



Parameter sharing: as kernel slides → able to capture same property in different image regions multiple feature detectors can be used to capture diff. properties

Convolution Layer



We do not have 1 kernel, we have \leq kernels

* $c = \text{nb. of channels}$ (kernels) in convolutional layer

notes terms padding: → added outside the image → allows us to use the kernel in all pixels of the image → prevent reduction in spatial dimensions

padding: extra border to input

size before kernel size padding

Output Size $O = \frac{W - K + 2P}{S} + 1$ stride

size after convolution

① padding: extra border to input added

② depth: nb of kernels (Filters) of Layer

③ Stride: step used to slide kernel on input

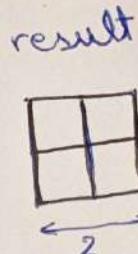
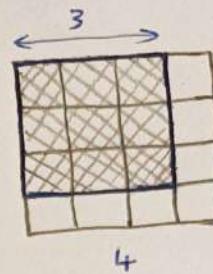
> 1 → down sampling input data
→ reduce spatial dimension output

CNN

Output size
after Convolution

$$O = \frac{W - K + 2P}{S}$$

stride = 1 $= \frac{4 - 3 + 0}{1} = 2$



Pooling: Further reduce dimensionality

Pooling Functions

- ① Average Pooling
- ② Max Pooling

2	5	1	5
1	8	3	2
0	6	1	2
5	3	3	0

max pooling

8	5
6	3

Avg pooling

.	.
.	.

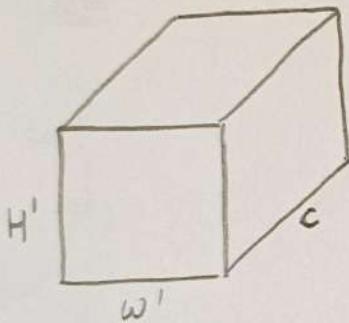
decimal averages

- * Suppose in a part of image we have special structure → kernel act with tolerance to its position

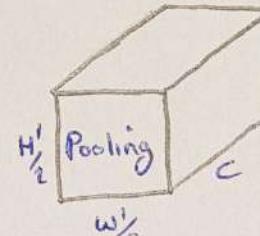


no parameters
to be
estimated

pooling: after convolution, use pooling functions
to aggregate elements of result grid (Max, Avg)
⇒ downsampling



convolution
+ ReLU



output of this layer
used as input to
next layer

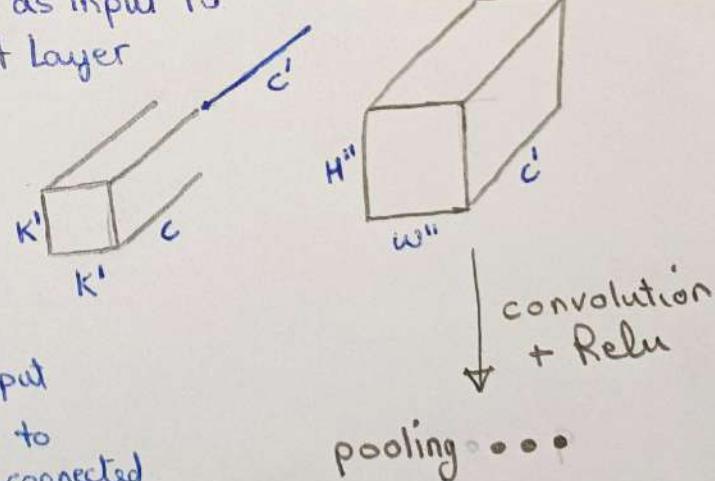
NB:

After CNN (feature recognition),
it is common to attach a DNN
for Classification.

- * Classification → need flattening

operation to convert output

of Conventional Layer to
proper input for Fully connected
Layer



pooling ...

- * $\text{Softmax}(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$ so output can be interpreted as probability

- * downsample image at different levels
 - ↳ low levels, kernel Learn Low features (small details)
 - ↳ high levels, kernel Learn high features (ex: faces)

MLE MAXIMUM LIKELIHOOD ESTIMATION

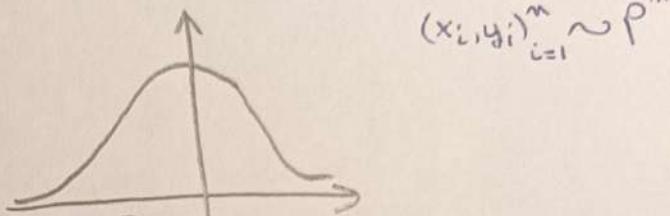
- * with Bayesian approaches, probability distribution are central objects of interest

minimization of Loss function: $\min \sum_{i=1}^n l(f_w(x_i), y_i) \approx \min \mathbb{E}[l(f_w(x), y)]$

$$\text{MLE: } \min e^{-\sum_{i=1}^n l(f_w(x_i), y_i)}$$

$$z_i = (x_i, y_i) \in \mathbb{R}^{d+1}$$

append



We assume P
to be gaussian
to reduce estimation
to mean

Gaussian

$$\star P = P_\theta \quad \theta \in \mathbb{R}^p \rightarrow \text{space of parameters}$$

$$\theta(\mu, \sigma)$$

$$\star P_\theta(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(z-\mu)^2}{2\sigma^2}}$$

search for θ that
maximize the likelihood

$$\star \text{if } \theta_1, \theta_2, \dots, \theta_m$$

$$(z_1, z_2, \dots, z_n) \sim P_\theta(z_1, z_2, \dots, z_n) = P_{\theta_1}(z_1) \cdot P_{\theta_2}(z_2) \cdots P_{\theta_m}(z_m)$$

Likelihood $\leftarrow \hat{L}(\theta) = \prod_{i=1}^n P_\theta(z_i)$

\prod : like Σ but for
products

$$\star L(\theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-\frac{(z_i-\mu)^2}{2\sigma^2}} \quad \text{where } \sigma = 1$$

$$= \frac{1}{(\sqrt{2\pi})^n} e^{-\frac{\sum_{i=1}^n (z_i-\mu)^2}{2}}$$

note

$$\log(ae^z) = \log a + \log e^z \\ = \log a + z$$

- * Instead of selecting θ & keep ones of highest L

optimize $\max \frac{1}{\sqrt{2\pi}} e^{-\frac{\sum_{i=1}^n (z_i-\theta)^2}{2}} \quad z \sim N(\theta, 1)$

take the Log:
 $\Rightarrow \max \log \hat{L}(\theta) = \text{constant} - \frac{\sum_{i=1}^n (z_i-\theta)^2}{2}$

$$\Rightarrow \max \sum_{i=1}^n (z_i - \theta)^2 / 2 \quad \therefore -\frac{d \log \hat{L}(\theta)}{d\theta} = 0$$

$$= \sum_{i=1}^n (z_i - \theta) = 0 \Rightarrow n\theta = \sum_{i=1}^n z_i$$

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n z_i$$

* normal distribution is characterized by mean μ & standard dev σ

note
mean: $\mu = \frac{1}{n} \sum_{i=1}^n x_i$

Pdf: probability dist. function: $f_{(\mu, \sigma)}(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

* Multivariate Gaussian N = normal distribution

$z \sim N(\theta, \sigma^2 I)$ * probability pdf of multivariate distr.

n set of gaussian dist. each with own mean

* $z \in N(\mu, \Sigma) \rightarrow P_\theta(z) = \frac{e^{-(x-\mu)^T \Sigma^{-1} (x-\mu)}}{c}$

$$\mu \in \mathbb{R}^d$$

$$\Sigma \in \mathbb{R}^{d \times d} = \mathbb{R}^{d^2}$$

$$\Theta(\mu, \Sigma) = \mathbb{R}^{d+d^2}$$

how variables are correlated? Σ = covariance Matrix is a symm matrix describing the covariances between diff. dimensions of multivariate distr

This pdf of multivariate dist. describe the likelihood of observing vector x given mean μ & covariance Σ

Multivariate Gaussian dist. \rightarrow describes joint distribution of multiple variables

μ is a mean vector

* case of n Gaussians, each with it's own μ & σ

$$P_\theta(z) = \frac{1}{c} \sum_{j=1}^n e^{-(x-\mu_j)^T \Sigma^{-1} (x-\mu_j)}$$

GMM: Gaussian Mixture Model.

Tails: of distribution, less probable values found at extremes far from mean \rightarrow heavy tails, light tails



not heavy tails \rightarrow tails decay rapidly \rightarrow extreme events less likely
 \rightarrow distribution concentration around mean

decay means value decrease as u move away from mean, values around mean more likely occur

MLE for Linear Regression

$$y_i = \omega^T x_i + \underbrace{\varepsilon_i}_{\text{error term}} \quad i=1 \dots n$$

MLE benefits → we know the probability distribution → compute the confidence of our prediction

* $P(y_i | x_i; \omega)$ $y_i \sim \mathcal{N}(\omega^T x_i, \sigma^2)$

$$= \mathcal{N}(\omega^T x_i, \sigma^2)$$

$$= c e^{-\frac{(y_i - \omega^T x_i)^2}{2\sigma^2}} \quad c = \text{constant}$$

* $P_w(y_1, \dots, y_n | x_1, \dots, x_n; \omega) = \prod_{i=1}^n c e^{-\frac{|y_i - \omega^T x_i|^2}{2\sigma^2}}$

Likelihood function in regression $\mathcal{L}(\omega) = c' e^{-\sum_{i=1}^n |y_i - \omega^T x_i|^2 / 2\sigma^2}$ Likelihood function for observing (x_i, y_i) given ω

Log $\mathcal{L}(\omega)$ = constant - expression expression is minimum, Likelihood increase

This is the term inside the exponent $\min \sum_{i=1}^n \frac{(y_i - \omega^T x_i)^2}{2\sigma^2}$

just Least square ERM

Maximize Likelihood equivalent to finding the ω that minimize the sum of square errors $\Rightarrow \max \mathcal{L}(\omega) \rightarrow \min \sum_{i=1}^n |y_i - \omega^T x_i|^2$
finding ω that maximize $\mathcal{L}(\omega)$

Linear Logistic Regression

$$\min \sum_{i=1}^n \log(1 + e^{-\omega^T x_i y_i}) \quad P_w(y_i | x_i)$$

MLE → does not make sense for classify

$$y_i \in \{1, -1\}$$

$$P(1|x) = \frac{1}{1 + e^{-\omega^T x}} \quad P(-1|x) = 1 - P(1|x) = \frac{e^{-\omega^T x}}{1 + e^{-\omega^T x}} = \frac{1}{1 + e^{\omega^T x}}$$

$$\hat{\mathcal{L}}(\omega) = \prod_{i=1}^n P_w(y_i) = \prod_{i=1}^n \frac{1}{1 + e^{-\omega^T x_i y_i}}$$

$$\log \hat{\mathcal{L}}(\omega) = \sum_{i=1}^n \log \frac{1}{1 + e^{-\omega^T x_i y_i}} = - \sum_{i=1}^n \log(1 + e^{-\omega^T x_i y_i})$$

Prior and Posterior

in bayesian θ is related to precision of prior $P(w)$

$$P(U|V) = \frac{P(V|U) P(U)}{P(V)}$$

↓
Posterior

prior = initial beliefs
posterior = updated beliefs
Likelihood = probability of observing data given parameters

Bayes theorem: $P(w|X_n, Y_n)$ represent the posterior dist. of coeff. w given X_n & Y_n

$$P(w|Y_n, X_n) = \frac{P(Y_n|X_n, w) P(w)}{P(Y_n)} \Leftrightarrow * P(Y_n|X_n, w) = e^{-\sum_{i=1}^n \frac{(y_i - w^T x_i)^2}{2\sigma^2}}$$

$$* P(w) = \frac{1}{C} e^{-\|w\|^2}$$

$$P(w|Y_n, X_n) = \frac{1}{C'} e^{-\sum_{i=1}^n \frac{(y_i - w^T x_i)^2}{2\sigma^2}} \cdot e^{-\|w\|^2}$$

$$= \frac{1}{C'} e^{-\left(\sum_{i=1}^n \frac{(y_i - w^T x_i)^2}{2\sigma^2} + \|w\|^2\right)}$$

full distribution of w

Maximum A Posterior

MAP

$$\max_{w \in \mathbb{R}^d} e^{-\left(\sum_{i=1}^n \frac{|y_i - w^T x_i|^2}{2\sigma^2} + \|w\|^2\right)}$$

map → parameter free

↑

$$\underline{\text{RLS}} \quad \min_{w \in \mathbb{R}^d} \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|^2$$

λ : param

* Ridge Regression can be derived as the Max. a Posteriori MAP for Linear regression

* Regularization Term in ridge regression determined by scaling Gaussian Prior

* Bayesian → allow computation of variance of parameter estimate → gives a more comprehensive view of parameter uncertainty

Sparsity

Linear model: $f_w(x) = w^T x = \sum_{j=1}^D w_j x_j$ $x \rightarrow (x^1, \dots, x^d)$

detect which variables are important for prediction ??

x^j : features of input x
measures --

which factors have determined our prediction

sparse: few of the coefficients are non-zero

$$\|w\|_0 = \{j \mid w_j \neq 0\}$$
 count non-zero components in w

l₀ norm

X_n : $n \times D$ matrix

$X^j \in \mathbb{R}^n$

$j=1 \dots D$ columns

$Y_n \in \mathbb{R}^n$ output vector

$$\text{square loss} \min_{w \in \mathbb{R}^D} \frac{1}{n} \sum_{i=1}^n (y_i - f_w(x_i))^2 + \lambda \|w\|_0$$

l₀ norm

→ approach is hard considering all training sets with all subsets of variables

We consider 2 approaches:

Greedy Method OMP orthogonal matching Pursuit

MP & Forward Stage-Wise Regression

- ① initialize $r_0 = Y_n$ initial residual
 $w_0 = 0$ initial coefficient vector
 $I_0 = \emptyset$ initial index set

- ② iterate $i = 1 \dots T-1$
and compute correlation coefficients: a_j

$$a_j = \frac{(r_{i-1}^T x^j)^2}{\|x^j\|^2} \quad \text{for}$$

- ③ Select K: $K = \operatorname{argmax}_{j=1 \dots D} a_j$

$$\text{note } v^i = \operatorname{argmin} \|r_{i-1} - x^i v\|^2$$

variable selection with: ① Larger projection on output
 ② best explains output vector in least square

④ Update Index: $I_i = I_{i-1} \cup \{K\}$ Union of previous iteration with K. K is index of best variable correlation

⑤ $w_i = w_{i-1} + \omega_K$ $\omega_K \overset{\text{K'th element}}{=} v_K e_K^T \neq \text{zero}$
 coefficient

Vector Update

w_K is a vector with only Kth component different than 0. zero

⑥ Residual Update:

$$r_i = r_{i-1} - X_K w_K$$

Orthogonal MP Same but

⑤ coefficient update

$$\underline{w_i} = \underset{w \in \mathbb{R}^D}{\operatorname{argmin}} \|Y_n - X_n M_i w\|^2$$

M_i

acts as selector

$M_i: D \times D$
 matrix

if j in set I

$$\Rightarrow (M_i w)^j = w^j$$

if solution is sparse & data matrix

has columns (features) not too correlated

\Rightarrow OMP recovers ^{with} high probability
 the right vector of coefficients

if j not in set I

$$(M_i w)^j = 0$$

LASSO and elastic net

LASSO : Least absolute shrinkage & Selection operator

based on convex relaxation : use ℓ_1 norm

$$\|w\|_1 = \sum_{j=1}^D |w_j|$$

another approach
to find sparse solution

Least Squares:

$$w_\lambda = \min \frac{1}{n} \sum_{i=1}^n (y_i - f_w(x_i))^2 + \lambda \|w\|_1$$

objective function
* convex but not differentiable
convex function: if any 2 points
in its domain, the line
connecting them lies above the graph of function

$\|w\|_1$ at zero ||
↓

Tools needed to find solution:

ISTA Iterative Soft Thresholding Algorithm

$$w_0 = 0$$

$$w_i = S_{\lambda\gamma} \left(w_{i-1} - \frac{2\gamma}{n} X^T (Y - Xw_{i-1}) \right)$$

↑
same
↓
Subdifferential
Soft thresholding

$$w_{t+1} = S_\gamma \left(w_t - \frac{2\gamma}{n} \hat{X}^T (\hat{X} w_t - \hat{Y}) \right)$$

iteration
 $i = 1 \dots T_{\max}$

Iteration run until convergence
is met:
* $\|w_i - w_{i-1}\| \leq \epsilon \rightarrow$ precise
→ max nb of iterations T_{\max}

w_i : coeff parameter vector at iteration i

γ : step size (learning rate)

n: nb of data points

λ : reg param

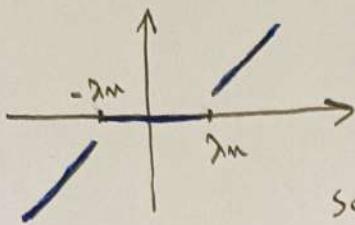
* The choice of λ has significant impact on resulting model

λ big → increases penalty & shrink more coefficients towards zero.
Less imp coefficients to zero (feature selection)

λ small → Less sparsity
→ Less regularization
→ potentially overfitting more prone to overfitting

$$\text{step size } \gamma = \frac{n}{2\|X^T X\|}$$

$$S_\lambda = \begin{cases} w_j - \lambda n & w_j \geq \lambda n \\ 0 & -\lambda n \leq w_j \leq \lambda n \\ w_j + \lambda n & w_j \leq -\lambda n \end{cases}$$



soft-thresholding is a component-wise operation on w

$$S_\alpha(u) = |u| - \alpha + \frac{u}{|u|}$$

$|u| - \alpha$: perform thresholding by
subs. α from each component of u

if $|u|$ greater than α
then $|u| - \alpha$
if $|u|$ less than α
then \rightarrow zero

$\frac{u}{|u|}$: ensure the sign of each component in result is same (preserve sign)

* $w_j > \lambda n$ means magnitude of coeff. w_j larger than λn \rightarrow soft-threshold operation result in w_j being shrunk to zero

* $-\lambda n < w_j < \lambda n$ magnitude between \rightarrow soft thresholding will shrink w_j toward zero but not become zero

* $w_j < -\lambda n$ - soft-th result w_j being pushed towards zero

Combination of these ² terms result in soft thresholding effect.

Tikhonov \rightarrow not sparse

LASSO \rightarrow not stable (might not be stable)

\Rightarrow hybrid algorithm between Tikhonov & Lasso

Elastic Net Algorithm

$$\min \frac{1}{n} \sum_{i=1}^n (y_i - f_w(x_i))^2 + \lambda (\alpha \|w\|_1 + (1-\alpha) \|w\|_2^2)$$

$\alpha \in [0, 1]$

Decision Trees

partition on the input Space

$$\left\{ A_j \subseteq X, j=1 \dots J \mid X = \bigcup_{j=1}^J A_j, A_j \cap A_i = \emptyset, i \neq j \right\}$$

each A_j
 is subset
 of input
 space X

 Union of all
 sets A_j equal
 entire input
 Space

 no elements in
 common

Dyadic Partition

- * each cell A_j , side length of each cube 2^{-k}
- * A_j is a cube * divide entire space into cubes.
- * partition A

partition based estimator learns a constant function

constant function: $\hat{c}_1, \dots, \hat{c}_J$ ^{respective} on each cell $A_1 \dots A_J$

$$\min_{c_1 \dots c_J} \sum_{j=1}^J \sum_{x_i \in A_j} (y_i - c_j)^2$$

We can compute each value of \hat{c}_j separately: $\min_{c_j \in R} \sum_{x_i \in A_j} (y_i - c_j)^2$
 → to find the constant value for each cell that minimize the sum of squared differences

$$j = 1, \dots, J \quad * \hat{c}_j = \arg \min_{c_j} \sum_{x_i \in A_j} (y_i - c_j)^2$$

and $\hat{c}_j = \frac{1}{n_j} \sum_{x_i \in A_j} y_i$ $n_j = \text{nb of points in } A_j \text{ cell}$

equivalently derived by piecewise constant function

$$\mathcal{H} = \left\{ f: X \rightarrow \mathbb{R} \mid f(x) = \sum_{j=1}^J c_j \mathbb{1}_{A_j}, \forall x \in X, c_1 \dots c_J \in \mathbb{R} \right\}$$

$\mathbb{1}_{A_j} = \begin{cases} 1 & x \text{ is in cell } A_j \\ 0 & \text{otherwise} \end{cases}$

$$\text{ERM: } \min_{\substack{f \in \mathcal{H} \\ \text{const}}} \sum_{i=1}^n (y_i - f(x_i))^2$$

$$\min_{c_1, \dots, c_j} \sum_{j=1}^J \sum_{x_i \in A_j} (y_i - c_j)^2$$



extend the idea beyond piecewise constant functions

$$\min_{\substack{f \in \mathcal{H} \\ \text{const}}} \sum_{i=1}^n (y_i - f(x_i))^2$$

• Piecewise Linear estimators

$$\mathcal{H}_{\text{lin}} = \left\{ f: X \rightarrow \mathbb{R} \mid f(x) = \sum_{j=1}^J (w_j^T x + b_j) \mathbb{1}_{A_j x}, \forall x, w_j \in \mathbb{R}^d, b_j \in \mathbb{R}, j=1 \dots J \right\}$$

$$\text{ERM of } \mathcal{H}_{\text{lin}} : \min \sum_{i=1}^n (y_i - f(x_i))^2$$

Computation:

solution require solving Least square in each cell

* idea Can now be extended:

$$\mathcal{H} = \left\{ f: X \rightarrow \mathbb{R} \mid f(x) = \sum_{j=1}^J g(x) \mathbb{1}_{A_j x}, \forall x \in X, g \in G \right\}$$

G : features, kernels, Neural network --

solution require solving a

$$\text{ERM on } \mathcal{H}_G : \min_{f \in \mathcal{H}_G} \sum_{i=1}^n l(y_i, f(x_i))$$

separate ERM in each cell (local erm in each cell)

* How to Select a partition?

Partition Tree

→ family of partitions

all space
 A_0, A_1, \dots, A_Q family of partitions
 $\xrightarrow{\text{increasing cardinality}}$ $A_0 = \{X\}$

$$\text{eg: } A = \{A_{1,2}; A_{3,1}; A_{3,2}; A_{3,3} \dots\}$$

each partition is a refinement of another partition

now for $A \in A_q \rightarrow c_1 \dots c_M$ s.t. $A = \bigcup_{j=1}^M c_j$ $J_q \leq M^q$

- * A_q is like a group of pieces (don't overlap) & together cover space
- * each partition is a refinement of another partition

Select / Construct partition from partition tree approach

① Uniform Partitioning: the depth parameter q control level of detail \rightarrow at each level of detail q , use cross validation to evaluate estimator performance \rightarrow finally select level of best performance

scale param. q can be seen as regularization parameter

② Adaptive Partitioning recursively building a partition \rightarrow refining a cell if the error in it is larger than threshold τ

let: $A \in A_q \rightarrow$

$\hat{L}_A = \min_{f \in H} \frac{1}{n_A} \sum_{x_i \in A} (y_i - f(x_i))^2$

number of input points in A

cells in A_q

So if: $\hat{L}_A \leq \tau$ less or equal threshold τ
 \Rightarrow we keep it in partition A kept in partition cell

Else $\hat{L}_A > \tau$ split A considering $c_1 \dots c_M$ smaller cells splitting
 and $c_1, c_2 \dots c_M \subset A_{q+1}$

Result partition contain cells of different size (refinement)

The next partition
 s.t. $A = \bigcup_{j=1}^M c_j$

Greedy approach in building decision Tree

$A(s, j)$: family of partitions made of 2 sets

$$A_1(j, s) \quad A_2(j, s) \quad j: j\text{-th component of } x$$

$$\{x \in X \mid x^j \leq s\} \quad \{x \in X \mid x^j > s\}$$

choose best parameters j and s that minimize

$$\min_{j,s} [L_{A_1(s,j)} + L_{A_2(s,j)}]$$

o get pair j_*, s_* solve above problem

we get/set $A_1 = A(s_*, j_*) \rightarrow$ family of partitions

made of 2 cells

$$\underline{A_{1,1}} \quad \underline{A_{1,2}}$$

Recursively apply the above scheme to $A_{1,1}$ & $A_{1,2}$

$$B_{2,1}(j, s) \quad B_{2,3}(j, s)$$

$$\{x \in A_{1,1} \mid x^j \leq s\} \quad \{x \in A_{1,2} \mid x^j \leq s\}$$

$$B_{2,2}(j, s) \quad B_{2,4}(j, s)$$

$$\{x \in A_{1,1} \mid x^j > s\} \quad \{x \in A_{1,2} \mid x^j > s\}$$

Clustering

Supervised Learning: $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$

$$x_i \in \mathbb{R}^d$$

$$y_i \in \mathbb{R}$$

Unsupervised Learning: $S = \{x_1, \dots, x_n\}$

* Clustering: grouping collection of objects into subsets/clusters

* assume the knowledge on the number of clusters $K \leq n$

$$k = C(x_i)$$

encoder C^*
achieve goal of
clustering

cluster label

$$K = \{1, \dots, K\}$$

* Encoder C^* is a function that assigns each x_i (data point) to a cluster of label K

assign similar or
close data points to same cluster

$$\text{Combinational Clustering} \rightarrow W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(i')=k} d(x_i, x_{i'})$$

distance or
dissimilarity

optimal Clustering $C^* = \arg \min W(C)$ unfeasible in practice
assignment

K-means

Cluster Centroids: $m_k \quad k=1, \dots, K$ ① initialize

$$\cup C_i = X$$

$$C_i \cap C_j = \emptyset \quad \forall i \neq j$$

$$C_i = \{x \in X / \|x - m_i\| \leq \|x - m_j\| \quad \forall j \neq i\}$$

$C(x_i) = \underset{\substack{\text{cluster} \\ \text{assignment}}}{\arg \min} \underset{1 \leq k \leq K}{\|x_i - m_k\|^2}$ ② assign data to centroids $c_1 \dots c_K$

$$m_1, \dots, m_K \quad \text{update}$$

③

go to ② and repeat until stop condition

Lloyd's algorithm

distance from point x to mean of C_i \leq than distance from x to other means

Note

$$m_k = \frac{1}{S_k} \sum x_i$$

calculate mean centroid of data points assigned to cluster, then update m_k

K-means aim to minimize the function:

$$J(C, m) = \sum_{k=1}^K \sum_{C(i)=k} \|x_i - m_k\|^2$$

* this is Kmeans \rightarrow iterative procedure \rightarrow ① assign data points to nearest centroid

- No guarantee to find global minimum
- algorithm can converge to local minimum

② update centroids

$K = \text{nb of clusters}$
 $C = \text{set of data points}$
 $m_k = \text{centroid of cluster } k$
 $x_i = \text{data point assigned to cluster } k$

Note: Quality of Solution
depend on initial config. of centroids

This is why K-means is run multiple times with different random initialization \rightarrow best solution of lowest cost selected

$$\min_{m_1, \dots, m_K} \sum_{i=1}^n \min_{j=1 \dots K} \|x_i - m_j\|^2$$

Update
Assign

K-means ++

Kmeans++ improves initialization step to enhance chances of converging to better solution

① try maximize distance between means (1st mean random)

② assign probability to each point points far from rest have higher prob. then do Lloyd's alg.

* first centroid is chosen as first point in permuted order (random selection from data points)

points chosen farther

\rightarrow Remaining centroids:

The probability dist. ensure that points with large squared distances (farther from existing centroids) are chosen

$$\operatorname{argmax}_{P(x)} P(x) = \frac{D(x)^2}{\sum_i D(x_i)^2}$$

pick the point with highest probability

- calculate distance from existing centroids to all data points. let $D(x)$ be shortest to nearest centroid
- calculate squared distances & compute probability for choosing new points as centers

\rightarrow repeat till all centroids chosen.