

Cornell Machine Learning Notes No. 1

- Machine Learning - Algorithms that improves themselves on some task with experience.
- Based on Statistics and Optimization, not logic.
 - Supervised learning, Unsupervised learning, reinforcement learning.

Terminology

Setup: $D = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2) \dots (\vec{x}_n, y_n)\} \subseteq \mathcal{X} \times \mathcal{Y}$

$(\vec{x}_i, y_i) \sim p$ → we never know.
i.i.d. ↗

$\mathcal{X} = \mathbb{R}^d$ classification
 $\rightarrow f(x_i) \rightarrow y_i$ ↗
regression.

Remember, the training data and test data must be drawn from the same distribution.



扫描全能王 创建

Cornell Machine Learning Notes No. 2

$D = \{ \dots \} \rightarrow \text{ML Algorithm} \rightarrow h \Rightarrow h(\vec{x}_i) \approx y_i$

$\vec{x}_{\text{new}} \rightarrow h \rightarrow y_{\text{new}}$

$h \in \mathcal{H} \rightarrow \text{hypothesis class.} \rightarrow \text{need to be chosen by human.}$

then we find the best h in the \mathcal{H} we choosed.
 how? \downarrow \rightarrow loss functions.

- 0/1-loss $l(h, D) = \frac{1}{n} \sum_{(x_i, y_i)} \delta[h(\vec{x}_i) \neq y_i]$
- squared loss $l(h, D) = \frac{1}{n} \sum (h(\vec{x}_i) - y_i)^2$
- absolute loss $l(h, D) = \frac{1}{n} \sum |h(\vec{x}_i) - y_i|$

Learning process : finding $h \in \mathcal{H}$ with small loss.

Generalization : $h \in \mathcal{H}$, s.t. $\forall (\vec{x}, y) \sim P, h(\vec{x}) \approx y$.

We want $E[l(h, \vec{x}, y)]_{(\vec{x}, y) \sim P}$ to be small.

but we don't have access to P , thus we use a

test set to estimate it, this E .



扫描全能王 创建

Cornell Machine Learning Note No. 3.

$D = \boxed{\text{train} \quad \text{valida...} \quad \text{test}}$ → can only be used to finally test the model, not to choose between models. or you will be overfitting the test set.
 From Law of Large Numbers.

$$l(h, D_{\text{test}}) = \frac{1}{n} \sum l(h(\vec{x}_i, y_i)) \xrightarrow{n \rightarrow \infty} E[l(h(\vec{x}, y))]$$

— How to choose h ? Is there a best h ?

- We have to make assumptions on the data, otherwise we can not learn any model of it.
 - Different h makes different assumptions.
 - Then we can generate a dataset that violates the assumptions and the model will perform arbitrarily bad on this set.
- \Rightarrow "No ~~Free~~ Lunch" theorem. there is no one best algorithm that is best at all datasets.

Note : The algorithms works best when their assumption are met.



扫描全能王 创建

K-nearest neighbors. : data points that are close together have the same labels.

For test point x .

$$S_x \subseteq D, |S_x| = k. \quad \text{if } (x', y') \in D - S_x, \text{ dist}(x, x') \geq \max_{x'' \in S_x} \text{dist}(x, x'')$$

(x'', y'') in S_x vote $\rightarrow y$.

the key is to use proper distance metric.

$$\text{dist}(x, z) = \left(\sum_{r=1}^d |[x_r]_r - [z_r]_r|^p \right)^{\frac{1}{p}}$$

$p=1 \rightarrow$ Manhattan

$p=2 \rightarrow$ Euclidean

$p \rightarrow \infty \rightarrow$ Maximum.

Curse of dimensionality.

In very high dimensional space, almost all the points will be on the edges and be far away from each other. There will not be such thing as "neighbor". One point might be a little close to the test point than another, but that means nothing.



Cornell Machine Learning Notes. No. 5.

k-nearest neighbor suffers a lot from the curse of dimensionality

but people are still using it to classify images. why?

The truth is, image data often lies in a small subspace of the seemingly high-dimensional space.

In low dimensional case, as n gets larger, k-nearest neighbor can be very very accurate, but that will take too long. [O(n^d)]

—
Perceptron \rightarrow the first machine learning algorithm.

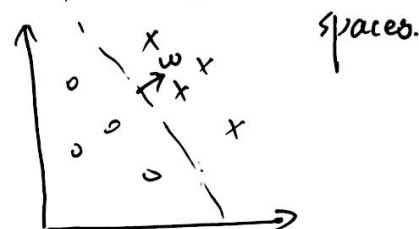
Assumption : the dataset is linearly separable.

\hookrightarrow almost always hold in high-dimensional space.

seldom hold in low-dimensional space.

Unlike k-nearest Neighbor, perceptron prefer high-dimensional

$$\mathcal{H} = \{x, w^T x + b = 0\}.$$



$$y = \text{sign}(w^T x + b) \quad \text{testing is fast for all possible points.}$$



扫描全能王 创建

Cornell Machine Learning Notes No. 6.

$\mathcal{Y} = \{-1, +1\}$. to learn w, b .

$$x_i \rightarrow \begin{pmatrix} x_i \\ 1 \end{pmatrix}, w \rightarrow \begin{pmatrix} w \\ b \end{pmatrix} \quad (w^T x_i' = w^T x_i + b)$$

$H = \left\{ x : w^T x = 0 \right\}$. to learn w .

Learning: $\vec{w} = \vec{0}$

while True :

$$m = 0$$

for every $(x, y) \in D$:

if $y w^T x \leq 0$:

$$w = w + yx$$

$$m = m + 1$$

if $m = 0$:

break.

END.

for $y = -1, w^T x \geq 0$

$$w = w - x$$

$$w^T x \rightarrow w^T x - x^T x < w^T x$$

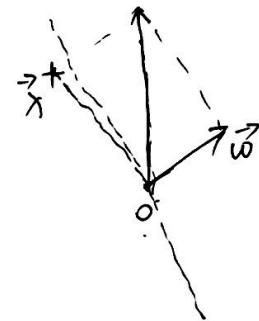
For $y = +1, w^T x \leq 0$

$$w = w + yx = w + x$$

$$w^T x \leq 0.$$

$$w^T x \rightarrow (w + x)^T x = w^T x + x^T x > w^T x$$

$$\dots \quad w^T x > 0$$



扫描全能王 创建

Cornell Machine Learning Notes No. 7.

$$\begin{aligned} y = +1 \Rightarrow w^T x > 0. \\ y = -1 \Rightarrow w^T x < 0 \end{aligned} \quad \left. \begin{array}{l} yw^T x > 0 \\ yw^T x < 0 \end{array} \right\} \text{or, update } w = w + yx.$$

Prove: this algorithm will end in finite steps.

$\exists w^*$. s.t. $\forall (x, y) \in D$, $yw^T x > 0$, we can assume $\|w^*\| = 1$

and $\|x\| \leq 1$ (since and we can rescale w and x without changing the sign of $yw^T x$)

High level intuition: if we can show that during iterations, $w^T w^*$ becomes more and more larger and larger, while $w^T w$ doesn't grow much, we can say that w is getting closer and closer to w^* , and they may eventually define the same hyperplane.

when update: $w \rightarrow w + yx \quad yw^T x \leq 0$.

$$w^T w^* \rightarrow (w + yx)^T w^* = w^T w^* + \underbrace{yw^* x}_{> 0} \geq w^T w^* + \gamma$$

$$\text{Margin: } \gamma = \min_{(x,y) \in D} |x^T w^*|$$

$$w^T w \rightarrow (w + yx)^T (w + yx) = w^T w + 2 \underbrace{y w^T x}_{\leq 0} + \underbrace{x^T x}_{\leq 1} \leq w^T w + 1$$

$$\text{After } M \text{ updates: } M\gamma \leq w^T w^* \leq \|w^*\| \|w^*\| = \sqrt{w^T w} \leq \sqrt{M} \Rightarrow M \leq \frac{1}{\gamma^2}$$



扫描全能王 创建

Cornell Machine Learning Notes. No. 8

$M \leq \frac{1}{\gamma^2}$, \Rightarrow datasets with large margin will converge faster.

If we can estimate $P(x, y)$, then we can get $P(y|x)$, then we can make predictions..

How to estimate it? Maximum Likelihood Estimation.

A unique, unfair coin, flip 10 times. got. 4 H + 6 T.

$$\text{MLE: } P(D; \theta) = \text{constant} \cdot \theta^4 \cdot (1-\theta)^6$$

$$P(H) \Rightarrow P(H) = \arg \max_{\theta} P(D; \theta)$$

$$\frac{\partial P(D; \theta)}{\partial \theta} = \arg \max_{\theta} \log P(D; \theta)$$

$$= \arg \max_{\theta} [4 \log \theta + 6 \log(1-\theta)]$$

$$\frac{\partial (4 \log \theta + 6 \log(1-\theta))}{\partial \theta} = \frac{4}{\theta} + \frac{6}{1-\theta} = 0 \Rightarrow \theta = \frac{4}{4+6}$$

$$\Rightarrow P(H) = \frac{n_H}{n_H + n_T}$$

This seems reasonable, but imagine we have little data and

they are all tails, we will get $P(H)_{\text{estimated}} = 0$, which is clearly not right.



扫描全能王 创建

Cornell Machine Learning Notes No. 9.

How to get around this? Smoothing.

If the dataset is relatively small, we don't trust it very much, instead we pretend that we have got n'_H of Heads and n'_T of tails before. Finally we will get $P(H) = \frac{n_H + n'_H}{n_H + n_T + n'_H + n'_T}$

Typically, n'_H, n'_T is set to be one.

This seems reasonable, there is a very nice way to deriving this. ✓
MAP

frequency view and bayesian view.

$$P(D; \theta)$$

$$P(D|\theta)$$

θ is a random variable.

which θ makes the data most possible?

θ is a parameter.

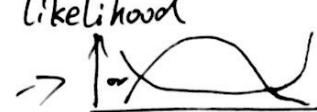
a $P(\theta)$ can be assumed.

this distribution encodes your prior belief about θ .

$$P(\theta|D) \rightarrow \text{likelihood} > P(\theta|D) \rightarrow \text{posterior}.$$

$P(\theta)$ → prior "Given the data, which is the most likely

$$P(\theta|D) = \frac{P(D|\theta) P(\theta)}{P(D)} = \text{constant. prior. likelihood}$$

gives value from 0 to 1. 

assume: $P(\theta) \sim \text{Beta distribution. } P(\theta) = \theta^{\alpha-1} (1-\theta)^{\beta-1}$

$$P(\theta|D) = \text{constant} \cdot P(\theta) \cdot P(D|\theta) = \theta^{\alpha-1} (1-\theta)^{\beta-1} \cdot \theta^{n_H} (1-\theta)^{n_T}$$

$$= \theta^{n_H + \alpha - 1} (1-\theta)^{n_T + \beta - 1}$$



扫描全能王 创建

Cornell Machine Learning Notes No. 10.

$$\text{Maximum } P(\theta | D) \Rightarrow P_{\text{MAP}} = \theta_{\text{MAP}} = \frac{n_H + \alpha - 1}{n_H + n_T + \alpha + \beta - 2}$$

\Rightarrow different approaches (smoothing and MAP) get the same result

From a bayesian point of view, by considering θ as a random variable, we can do sth amazing:

almost always unsolvable. ($\int_{\theta} p(y|\theta) p(\theta|D) d\theta$)

$$P(y=\text{HEAD} | D) = \int_{\theta} p(\text{heads}, \theta | D) d\theta = \int_{\theta} p(\text{heads}, | \theta, D) \cdot p(\theta | D) d\theta.$$

This is to say, by seeing θ as a random variable, we can take into account all the possible models and make the final prediction.

To estimate $P(x, y)$. : MLE $\theta = \arg \max_{\theta} P_{\theta}(D)$

MAP $\theta = \arg \max_{\theta} P(\theta | D)$

A more complex case of MLE.

$P(x)$. $P(D \sim P(y|x))$ $x \rightarrow$ a single feature

$$P(y|x) = \frac{\sum_{i=1}^n I(x_i = x) \cdot I(y_i = y)}{\sum_{i=1}^n I(x_i = x)} \left(\begin{array}{c} \frac{P(Y=y, X=x)}{P(X=x)} \\ \vdots \end{array} \right)$$

If \vec{x} is a ~~is~~ d -dimensional vector?

It is unrealistic to see exactly the same \vec{x}



扫描全能王 创建

↳ to rescue \rightarrow Naive Bayes : $P(y=y|x=x) = \frac{P(x=x|y=y) \cdot P(y)}{P(x=x)}$

$P(x=x)$ \rightarrow Normalizer. $P(x=y)$. \rightarrow ratio of y . in all y .

$P(x=x|y=y)$ assume the features are independent given the label. \rightarrow this is where the word "Naive"

It turns out in use, this assumption doesn't come from break things too much

$$h(\vec{x}) = \arg \max_y P(y|\vec{x}) = \arg \max_y P(x|y) P(y)$$

$$= \arg \max_y \left[\prod_{d=1}^D P(x_d|y) \right] P(y)$$

$$= \arg \max_y \left[\log P(y) + \sum_{d=1}^D \log(P(x_d|y)) \right]$$

A few cases of Naive Bayes.

$$\bullet x_d \in \{v_1, v_2, \dots, v_k\} \quad P(x_d=j|y=c) = [\theta_{j|c}]_d \quad \sum_{j=1}^k [\theta_{j|c}]_d = 1$$

\downarrow
k categories.

$$[\theta_{j|c}]_d = \frac{\sum_{i=1}^n I(y=c) I(x_d=v_j)_d}{\sum_{i=1}^n I(y=c) + k_d}$$

$$\bullet x_d \in \{0, 1, 2, \dots, m\} \quad m = \sum_{d=1}^D x_d$$

$$P(\vec{x}|m, y=c) = \frac{m!}{x_1! \cdot x_2! \cdots x_D!} \prod_{d=1}^D \theta_{x_d|c}^{x_d}$$



Logistic Regression.

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

- Naive Bayes \rightarrow estimate $P(x|y)$, then $P(y|x)$

- Logistic Regression \rightarrow estimate $P(y|x)$

$$\text{MLE: } \underset{w, b}{\operatorname{argmax}} \prod_{i=1}^n P(y_i | \vec{x}_i; w, b) = \underset{w, b}{\operatorname{argmax}} \sum_{i=1}^n \log P_{w, b}(y_i | \vec{x}_i; w, b)$$

$$= \underset{w}{\operatorname{argmax}} \sum_{i=1}^n \log P_w(y_i | \vec{x}_i, w) = \underset{w}{\operatorname{argmax}} \log \frac{1}{1+e^{-y^T w}}$$

$$= \underset{w}{\operatorname{argmax}} - \sum_{i=1}^n (1 + e^{-y^T w}) = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n (1 + e^{-y^T w})$$

$$\text{MAP: } \underset{w}{\operatorname{argmax}} P(D|w) \cdot P(w) \quad P(w) \sim N(\vec{0}, \lambda I)$$

$$= \underset{w}{\operatorname{argmax}} \left(\prod_{i=1}^n \frac{1}{1+e^{-y^T w}} \right) \cdot e^{-\frac{w^T w}{2\lambda}}$$

$$= \underset{w}{\operatorname{argmax}} \left[\sum_{i=1}^n \log (1 + e^{-y^T w}) - \frac{w^T w}{2\lambda} \right]$$

$$= \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log (1 + e^{-y^T w}) + \alpha w^T w$$



$\min L(w)$

$$(w + \vec{s}) \approx (\vec{w}) + g(\vec{w})^T \vec{s} + \frac{1}{2} \vec{s}^T H(\vec{w}) \vec{s} + \dots$$

only holds when \vec{s} is really small step.

Gradient descent. $\vec{s} = -\alpha \vec{g}(w)$

$$(w + \vec{s}) \approx (\vec{w}) - \alpha \vec{g}(w)^T \vec{g}(w) \rightarrow \text{smaller}$$

Newton's Method. \rightarrow works better when already near the minimum.

$$(w + \vec{s}) \approx (\vec{w}) + g(\vec{w})^T \vec{s} + \frac{1}{2} \vec{s}^T H(\vec{w}) \vec{s} + \dots$$

$$\frac{\partial L(w + \vec{s})}{\partial \vec{s}} = g(\vec{w}) + H(\vec{w}) \vec{s} = 0 \Rightarrow \vec{s} = -H(\vec{w})^{-1} g(\vec{w})$$

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

Linear Regression. : $y = w^T x + \epsilon$. $\epsilon \sim N(0, \sigma^2)$

$$P(y_i | \vec{x}_i; w) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(w^T x_i - y_i)^2}{2\sigma^2}}$$

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

MLE: $\arg \max_w \prod_{i=1}^n P(y_i | \vec{x}_i; w) = \arg \max_w \prod_{i=1}^n e^{-\frac{(w^T x_i - y_i)^2}{2\sigma^2}}$

$$= \arg \max_w \sum_{i=1}^n -(w^T x_i - y_i)^2 = \arg \min_w \sum_{i=1}^n (w^T x_i - y_i)^2$$



MAP: $P(w|D) = \text{Constant} \cdot P(D|w) \cdot P(w)$

$$\underset{w}{\operatorname{argmax}} \left[\prod_{i=1}^n P(y_i | \vec{x}_i, w) \right] \cdot P(w) = \underset{w}{\operatorname{argmax}} \left(\prod_{i=1}^n e^{-\frac{(w^\top \vec{x}_i - y_i)^2}{2\sigma^2}} \right) \cdot e^{\frac{w^\top w}{2}}$$

$$= \underset{w}{\operatorname{argmax}} \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (w^\top \vec{x}_i - y_i)^2 - \frac{1}{2} w^\top w \right)$$

$$= \underset{w}{\operatorname{argmin}} \sum_{i=1}^n (w^\top \vec{x}_i - y_i)^2 + \lambda w^\top w$$

- Support Vector Machines. - Maximum Margin

Margin: the closest distance between the points and hyperplane.

$$\mathcal{H} = \{w^\top x + b = 0\}$$

$$\gamma = \min_{(x,y) \in D} \frac{|w^\top x + b|}{\sqrt{w^\top w}}$$

$$\text{Margin}(w, b) = \max_{w,b} \gamma(w, b)$$

$$w^\top x_p + b = 0 = w^\top (\vec{x} - \vec{d}) + b$$

$$w^\top (\vec{x} - 2\vec{w}) + b = 0$$

$$\Rightarrow \alpha = \frac{w^\top x + b}{w^\top w}$$

$$\Rightarrow d = \sqrt{d^\top d} = \frac{|w^\top x + b|}{\sqrt{w^\top w}}$$

$$\max_{w,b} \gamma(w, b) \quad \text{s.t. } \forall i, y_i(w^\top x_i + b) \geq 0$$

$$\max_{w,b} \left[\min_{(x,y) \in D} \frac{|w^\top x + b|}{\sqrt{w^\top w}} \right] \quad \text{s.t. } \forall i, y_i(w^\top x_i + b) \geq 0$$

$$\max_{w,b} \frac{1}{\sqrt{w^\top w}} \left[\min_D |w^\top x + b| \right] \quad \text{s.t. } \forall i, y_i(w^\top x_i + b) \geq 0$$



Cornel Machine Learning Notes. No. 15

Fix $\min_{\mathbf{w}, b} |\mathbf{w}^T \mathbf{x}_i + b| = 1$. (can be done by rescaling of \mathbf{w} and b)

$$\max_{\mathbf{w}, b} \frac{1}{\sqrt{\mathbf{w}^T \mathbf{w}}} \text{ s.t. } \rightarrow \min_{\mathbf{w}, b} \mathbf{w}^T \mathbf{w} \text{ s.t. } \begin{cases} y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 0 \\ \min_{\mathbf{w}, b} |\mathbf{w}^T \mathbf{x}_i + b| = 1 \end{cases}$$

$$\Leftrightarrow \min_{\mathbf{w}, b} \mathbf{w}^T \mathbf{w} \text{ s.t. } \forall i, y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

What if the dataset is not linearly separable?

$$\min_{\mathbf{w}, b} \mathbf{w}^T \mathbf{w} \text{ s.t. } \forall i, y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i, \xi_i \geq 0 \quad \text{hinge loss.}$$

$$+ C \sum_i \xi_i \quad \xi_i = \begin{cases} 0 & \text{if } y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \\ 1 - y_i (\mathbf{w}^T \mathbf{x}_i + b) & \text{if } y_i (\mathbf{w}^T \mathbf{x}_i + b) < 1 \end{cases}$$

$$\min_{\mathbf{w}, b} \mathbf{w}^T \mathbf{w} + C \sum_i \max(1 - y_i (\mathbf{w}^T \mathbf{x}_i + b), 0)$$

regularizer.

loss.

— Objective = $\frac{1}{n} \sum_{i=1}^n \text{Loss}(\mathbf{w}) + \lambda \text{Regularizer}(\mathbf{w})$ (apply to many models)

— Loss functions:

Hinge Loss. $\max[1 - h(\mathbf{x}_i) y_i, 0]^p$ $p=1 \rightarrow \text{SVM}$

Log Loss $\log(1 + e^{-y_i h(\mathbf{x}_i)})$ Logistic regression.



扫描全能王 创建

Cornell Machine Learning Notes. No. 16.

Exp Loss	$e^{-y_i h(x_i)}$	Adaboost
0/1 Loss	$\delta(h(x_i) \neq y_i)$	

Regression. $y \in \mathbb{R}$

Squared Loss - $(h(x_i) - y_i)^2$ - estimates mean.

absolute Loss - $|h(x_i) - y_i|$ - estimate median.

Juher Loss - $\begin{cases} \frac{1}{2}(h(x_i) - y_i)^2 & \text{if } |h(x_i) - y_i| < \delta \\ \delta(|h(x_i) - y_i| - \frac{\delta}{2}) & \text{otherwise} \end{cases}$

log-cook. Loss - ~~$\log \frac{e^{x_i} + e^{-x_i}}{2}$~~ $\log \left(\frac{e^{h(x_i) - y_i} + e^{-(h(x_i) + y_i)}}{2} \right)$.

Regularization.

$$w^T w$$

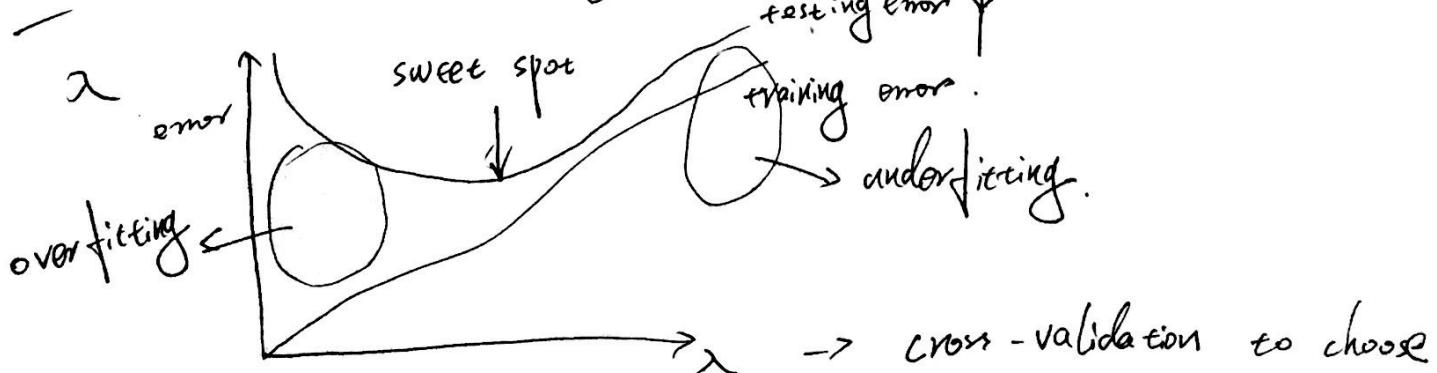
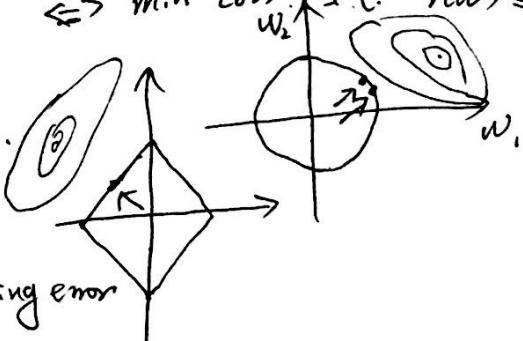
$$\min_{w^T w} \text{loss} + \lambda r(w) \Leftrightarrow \min_{w^T w} \text{loss} \text{ s.t. } r(w) \leq B.$$

$$r(w) = w^T w$$

L_2 -regularization.

$$r(w) = \|w\|$$

L_1 -regularization



扫描全能王 创建

Cornell Machine Learning Notes [lo. +]

Neural Networks.

$h(x) = \underline{w^T \phi(x)}$ $\phi(x) \rightarrow$ to be learned
 $h(x) = \phi^T w + b$ $x \rightarrow \phi(x)$ $\phi(x) = \sigma(Ax + b)$ some non-linear function.
mapping x into some other space and learn that mapping.

$$h(x) = w^T \phi(x) + b \quad \phi(x) = \sigma(ux + c)$$

$$h(x) = w^T \phi(x) + b \quad \phi(x) = \sigma(u\phi(x) + c) \quad \phi(x) = \sigma(\dots)$$



扫描全能王 创建

Cornell Machine Learning Note No. 18.

Boosting

$$H(x) = \sum_{t=1}^T 2_t h_t(x)$$

$$y = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix} \quad (x_1, x_2, x_3) \rightarrow \begin{pmatrix} H(x_1) \\ H(x_2) \\ H(x_3) \end{pmatrix} \quad \begin{pmatrix} h(x_1) \\ h(x_2) \\ h(x_3) \end{pmatrix}$$

Gradient Descent in the Function Space.

$$H(x) = \sum_{t=1}^T 2_t h_t(x)$$

$$h = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \quad l(H + 2h) \approx \underset{h \in \mathcal{H}}{\operatorname{argmin}} \left[l(H) + \left\langle \frac{\partial l}{\partial H}, 2h \right\rangle \right]$$

$$= \underset{h \in \mathcal{H}}{\operatorname{argmin}} \quad \left\langle \frac{\partial l}{\partial H}, 2h \right\rangle = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \quad \sum_{i=1}^n \frac{\partial l}{\partial H(x_i)} \cdot 2h(x_i)$$

$$\text{if } l(H) = \frac{1}{n} \sum_{i=1}^n (H(x_i) - y_i)^2 \quad \frac{\partial l}{\partial H(x)} = \begin{pmatrix} H(x_1) - y_1 \\ H(x_2) - y_2 \\ \vdots \\ H(x_n) - y_n \end{pmatrix}$$

$$\frac{\partial l}{\partial H(x)} = \frac{1}{n} \cdot 2(H(x) - y)$$

$$h \text{ in the direction of } y - H(x) = H(x) - y$$



扫描全能王 创建

Cornell Machine Learning Note No. 19.

Gradient boosting trees

of GBR, $h \in \text{CART}$ of limited depth, 2 constants

$$\arg\min h \quad 2 \sum_i r_i h(x_i)$$

$$r_i := \frac{\partial L}{\partial h(x_i)} \quad t_i = -r_i$$

$$\arg\min_h - \sum_i t_i h(x_i)$$

$\sum_i t_i^2$ has nothing to do with h .

$$\sum_i h(x_i)^2 = C$$

$$\arg\min_h \sum_i h(x_i)^2 - 2 \sum_i t_i h(x_i) + \sum_i t_i^2$$

$$\arg\min_h \sum_i [h(x_i) - t_i]^2 \rightarrow \text{still squared loss.}$$

$$\hat{y}(x) = \sum_{t=1}^T 2_t h_t(x)$$



扫描全能王 创建