

11 / 09 / 23

- Machine Learning

- Association

- Supervised Learning

- Classification » Labels : $\{1, 2, 3, \dots\}$ ^{categorical}
 - Regression » Labels : $f(x)$
Targets

- Unsupervised Learning » Grouping / Sorting

- Reinforcement Learning

a) Supervised (Regression)

b) Unsupervised

c) Speaker \rightarrow (Classification)

d) Supervised
Reinforcement

Machine Learning

Supervised Learning

Regression
classification

Training set $X = \{\hat{x}^t, r^t\}_{t=1}^N \sim$ total samples

$$\hat{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim 2 \text{ features} \quad r = \begin{cases} 1 & \text{if } \hat{x} \in C \\ 0 & \text{else} \end{cases}$$

Let $x_1 = \text{price}$ & $x_2 = \text{engine power}$,
 then, class C annotates training
 samples to a class by:

$$\therefore (p_1 \leq \text{price} \leq p_2) \wedge (e_1 \leq \text{engine power} \leq e_2)$$

- input ~ algorithm ~ logic ~ hypothesis
 $h(x) \leftarrow$

$$h(x) \triangleq \begin{cases} 1 & \text{if } h \text{ says } x \text{ is +ive} \\ 0 & \text{if } h \text{ says } x \text{ is -ive} \end{cases}$$

- error of h on H :

$$E(h | X) = \sum_{t=1}^N \mathbb{1}(h(\hat{x}^t) \neq r^t)$$

most
→ specific

- $h \in H$: b/w
 - \mathcal{G}_1 makes up version space
 - general
most

Margin

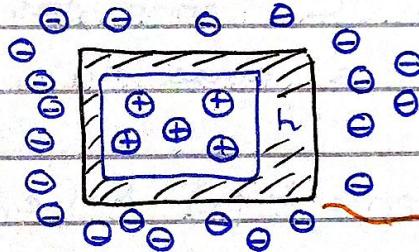
Choose h with largest margin
mean of
 S and G

Most { Specific General } Hypothesis

Probably Approximately Learning Correct (PAC)

How many sample
 N should we have
such that h has
error at most ϵ ;

with probability $1 - \delta$ $\rightarrow \delta \leq 1/2$
 $\epsilon > 0$



error with
respect to
hypothesis h
(not class)

- Each side (strip) is at most $\epsilon/4$
- P_r that we miss a strip = $1 - \epsilon/4$
- P_r that N instances miss a strip = $(1 - \epsilon/4)^N$
- P_r that N instances miss 4 strips = $4(1 - \epsilon/4)^N$

$$\left[4(1 - \epsilon/4)^N \leq \delta, \text{ solve for } N \right]$$
$$N \geq (4/\epsilon) \log(4/\delta)$$

Machine Learning

Noise and Model Complexity

→ simple models are usually better to use than a highly complex one

↳ irrelevant features to task

; such as, hair style for voice recognition

outliers generate - [bias when true outlier in classed (but)]
 error when noise

- if data lies outside predicted class box, reject it; do not give wrong information
- if data lies inside two overlapping clusters, reject it as well

Multiclass Labels

$$x = \{\hat{x}^t, r^t\}_{t=1}^N$$

$$r_i^t = \begin{cases} 1 & \text{if } \hat{x}^t \in C_i \\ 0 & \text{if } \hat{x}^t \in C_j, j \neq i \end{cases}$$

C ₁	0	0	0
C ₂	0	0	1
C ₃	0	1	0
:	:		:	
C _K	1	0	0

$$h_i(x^t) = \begin{cases} 1 & \text{if } \hat{x}^t \in C_i \\ 0 & \text{if } \hat{x}^t \in C_j, j \neq i \end{cases}$$

Regression

$$r^t = f(x^t) + \epsilon$$

$$\text{Linear Regression} \rightsquigarrow f(x^t) = \underbrace{w_1 x^t + w_0}_{\text{Model learns}}$$

$$\text{MSE} \triangleq E(f|X) = \frac{1}{N} \sum_{t=1}^N [r^t - g(x^t)]^2$$

Mean Squared Error

w₁ and w₀
for best fit

Triple Trade-off

- [complexity of H is $C(H)$]
- [Training Set Size, N]
- [Generalization Error, E , on new data]

- Test Data $G \setminus D \triangleq \{ \text{Training Data} \cup \text{Validation Data} \}$

Practice

- We use mean (average) of S and G_1 to minimize generalization error E .
- Yes, we can. Circle will have parameters {center, r }. Ellipse has {center, major axis, minor axis}.
- Outliers yield high error (\sim variance formula); for robust regression, we opt for Huber / Pseudo-Huber loss which integrates a user specified parameter s , and on its basis, we have L_1 loss for outliers and L_2 loss for inliers.

Sir : $f(x) = [1/\sqrt{2\pi}] e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2} \rightarrow \text{Gaussian Dist.}$

Assume distribution is Gaussian;

MSE serves us well

25/09/23

Machine Learning

Bayesian Decision Theory

Lack of info about origin of data compels us to believe that the problem is best modeled as statistical processes.
stochastic

Coin Toss

[observable variable $\hat{=}$ outcome
unobservable variable $\hat{=}$ temperature, humidity, etc.

$\rightsquigarrow z$: unobserved x : observable

$x = f(z)$; $f(\cdot)$ is a deterministic function

We don't have $f(\cdot)$ so we use probability $P(x = x)$

Tossing a coin $\in \{H, T\}$

RV $x \in \{0, 1\}$

Bernoulli $P(x = 1) = p_0^x (1 - p_0)^{1-x}$

\rightarrow Finding the approximator p_0 as \hat{p}_0

Estimation

$$\hat{p}_0 = \frac{\#\{H\}}{\#\{N\}}$$

no. of heads ↗ ↘ total tosses

$$= \frac{\sum x^t}{N}$$

Next toss is heads if $\hat{p}_0 > 1/2$

Bayes' Rule

$$P(C|x) = \frac{P(C) P(x|C)}{P(x)}$$

prior ↗ likelihood ↘
posterior ↗ evidence ↘

Prior $P(C=0) + P(C=1) = 1$

Evidence $P(x) = P(x|C=1)P(C=1) + P(x|C=0)P(C=0)$

Posterior $P(C=0|x) + P(C=1|x) = 1$

Likelihood $P(x|C=0) + P(x|C=1) \neq 1$

Bayes incorporates \rightarrow initial belief

and evidence to

make revised decisions (posterior)

Example 1 : $C = 1$; innocent

$C = 0$; guilty

\hat{x} = blood matches

$$P(C=1) = 0.4 \quad P(C=0) = 0.6$$

$$P(x|C=0) = 0.95 \quad P(x) = \text{Formula} = 0.61$$

$$P(x|C=1) = 0.1$$

$$P(C=1|x) = (0.4)(0.1)/0.61 = 0.0655$$

Example 2 : $C=1$ +ive $x = \text{PCR +ive}$
 $C=0$ -ive

$$P(C=1) = 0.02 \quad P(C=0) = 0.98$$

$$P(x|C=1) = 0.85$$

$$P(x|C=0) = 0.05$$

$$P(x) = \text{Formula} = 0.066$$

$$P(C=1|x) = 0.257$$

Bayes Rule $K \geq 2$

$$P(C_i|x) = \frac{P(x|C_i)P(C_i)}{P(x)}$$

$$= \frac{P(x|C_i)P(C_i)}{\sum_{k=1}^K P(x|C_k)P(C_k)}$$

$$\sum_{k=1}^K P(C_k) = 1$$

Choose C_i if $P(C_i|x) = \max_k P(C_k|x)$

Machine Learning

Losses and Risks

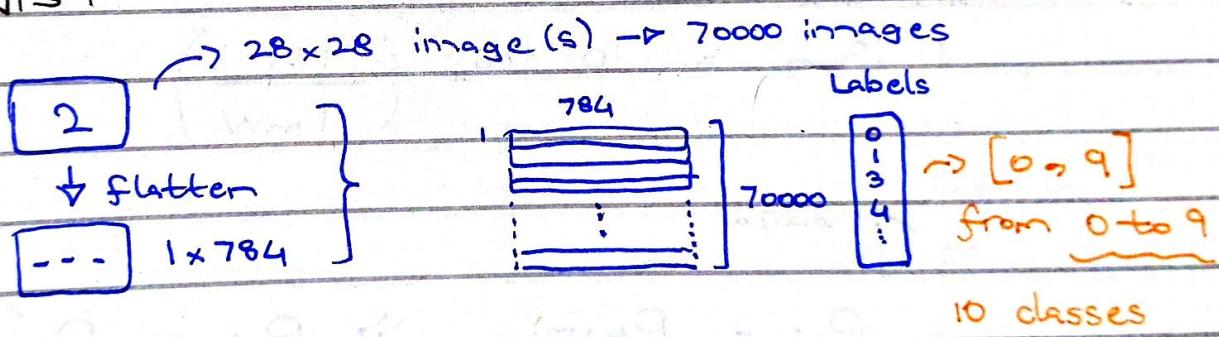
Actions \hat{a}_i

Loss of a_i when state is $C_k \hat{x}_{ik}$

Expected Risk $R(a_i|x) = \sum_{k=1}^K \lambda_{ik} P(C_k|x)$

choose a_i if $R(a_i|x) = \min_k R(a_k|x)$

MNIST



Training : 60000 samples Test : 10000 samples

- Suppose 5000 labels of class '0' exist, then;

$$p(C=0) = \frac{5000}{60000} = 0.083$$

↳ in training

- Normalize the data [0 - 1] to map to pdf better; and to get the data to a same scale.

Machine Learning

$$p(x|c) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2} \frac{(x-\mu_c)^2}{\sigma^2}\right)$$

$$\text{model}_c \sim [\{\mu_c, \sigma_c\}, p(c)]$$

$|\Sigma|^2 \rightarrow$ square of determinant of covariance matrix \rightarrow replaces σ



multivariate

Model \sim model_c

$$p(x|c) = \frac{1}{\sqrt{2\pi}|\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x-\mu) \Sigma^{-1} (x-\mu)^T\right)$$

verify formula from web $\rightarrow d$ is shape of Σ ($d \times d$)

$$\text{Cov} = \Sigma = \sum_i [(x - \mu_c^{(i)}) (x - \mu_c^{(i)})^T] = 784 \times 784$$

$$\mu_c = 784 \times 1$$

• Association Rules

$$X \rightarrow Y$$

$\left\{ \begin{array}{l} \text{implies association and not} \\ \text{necessarily causation} \end{array} \right\}$

→ Measures

$$\left. \begin{array}{l} \text{Support} \triangleq P(X,Y) = (X \text{ and } Y) / \text{Number} \\ \text{Confidence} \triangleq P(Y|X) = (X \text{ and } Y) / X \\ \text{Lift} \triangleq \frac{P(Y|X)}{P(Y)} \text{ or } \frac{P(X,Y)}{P(X)P(Y)} \end{array} \right\}$$

For cases such as, $X \rightarrow Y, Z$ and $X, Y \rightarrow Z$, we get;

$$c(X \rightarrow Y, Z) = \underbrace{X \text{ and } Y \text{ and } Z}_{\times} \quad \left. \begin{array}{l} \text{Do from} \\ \text{next wherever} \\ \text{written} \end{array} \right\}$$

$$c(\underbrace{X, Y \rightarrow Z}) = \underbrace{X \text{ and } Y \text{ and } Z}_{\times \text{ and } Y} \rightarrow \text{Not confirmed}$$

all events are independent of each other

oo | Search it up

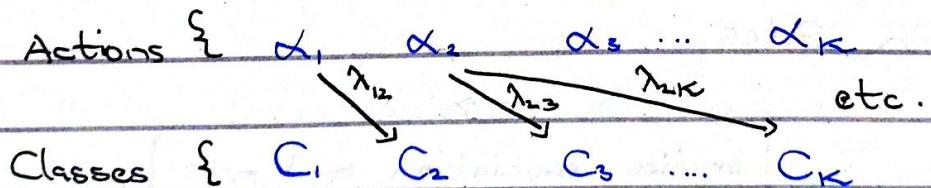
Risks & Losses : Continued

$$\text{Let } \lambda_{ik} = \begin{cases} 0 & \text{if } i \neq k \\ \lambda & \text{if } i = k+1 \\ 1 & \text{if } i = k \end{cases} \quad R(\alpha_{i+1} | x) = \lambda \quad R(\alpha_i | x) = 1 - P(C_i | x)$$

$$\rightarrow \lambda_{ik} \rightarrow \{\lambda_{01}, \lambda_{00}, \lambda_{10}, \text{etc.}\}$$

$$\text{where } i = 1, 2, 3, \dots, K$$

$$k = 1, 2, 3, \dots, K$$

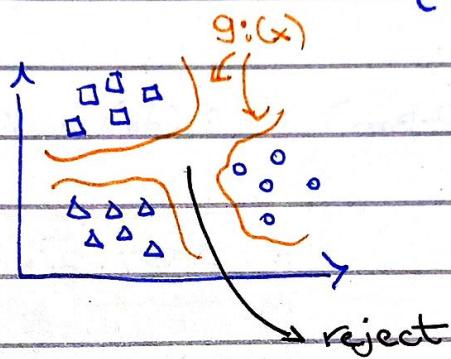


$$\rightarrow \text{In 0/1 loss ; } \lambda_{ik} = 0 \text{ if } i = k, 1 \text{ otherwise}$$

Machine Learning

• Discriminant Functions

$$i = 1, 2, \dots, K \gg g_i(x) = \begin{cases} -R(\alpha_i | x) \\ P(C_i | x) \\ P(x | C_i) P(C_i) \end{cases}$$



$$R_i = \{x \mid g_i(x) = \max_k g_k(x)\}$$

Decision Regions

• Dichotomizer $K = 2$

Polychotomizer $\{K \geq 3\}$

$$g(x) = g_1(x) - g_2(x)$$

Choose $\begin{cases} C_1 & \text{if } g(x) > 0 \\ C_2 & \text{otherwise} \end{cases}$

→ Log Odds : $\log \left[\frac{(P(C_1 | x))}{(P(C_2 | x))} \right]$

"gib me time"

buh
i-i

Association Measures

Multiple Variables

$$c(x, y \rightarrow z) = \frac{P(x, y, z)}{P(x) P(y)}$$

$$l(x, y, \rightarrow z) = \frac{P(x, y, z)}{P(x) P(y) P(z)}$$

According
to
Sir

Exercise 2

$C=0$ No Disease

$C=1$ Disease

x = +ive test

$$P(C=1) = 1/10^6 \quad P(C=0) = 0.\overline{99}$$

$$P(x|C=1) = 0.99$$

$$P(x|C=0) = 1/10^3$$

$$\Rightarrow P(x) = P(C=1)P(x|C=1) + P(C=0)P(x|C=0)$$
$$= 1e-3 \quad \text{Verify}$$

$$\Rightarrow P(C=1|x) = \frac{P(C=1) P(x|C=1)}{P(x)}$$
$$= 9.89 e-4 \quad \text{Verify}$$

Exercise 3

$$\text{Given: } \log \left[\frac{P(C_1|x)}{P(C_2|x)} \right] \gg \log \left[\frac{P(x|C_1) P(C_1)}{P(x|C_2) P(C_2)} \right]$$

$$g(x) = \log \left[\frac{P(x|C_1)}{P(x|C_2)} \right] + \log \left[\frac{P(C_1)}{P(C_2)} \right]$$

Exercise 4

$$\lambda_{11} = \lambda_{22} = 0 ; \quad \lambda_{12} = 10 ; \quad \lambda_{21} = 5$$

Write the optimal decision rule.

	C ₁	C ₂
α ₁	0	10
α ₂	5	0

$$\rightarrow \text{Risk } R(\alpha_i|x) = \sum_{k=1}^K \lambda_{ik} P(C_k|x)$$

$$R(\alpha_1|x) = \lambda_{11} P(C_1|x) + \lambda_{12} P(C_2|x) \\ = 0 + 10 P(C_2|x)$$

$$R(\alpha_2|x) = \lambda_{21} P(C_1|x) + \lambda_{22} P(C_2|x) \\ = 5 P(C_1|x) + 0$$

$$\rightarrow \text{Known: } P(C_2|x) = 1 - P(C_1|x)$$

$$R(\alpha_1|x) = 10 [1 - P(C_1|x)]$$

- Choose C_1 if $R(\alpha_1 | x) < R(\alpha_2 | x)$

$$\rightarrow 10 - 10P(C_1 | x) < 5P(C_1 | x)$$

$$\rightarrow \underline{P(C_1 | x) > \frac{2}{3}}$$

Univariate and Multivariate Analysis

Parameter Estimation

Mean and variance are enough parameters to estimate an underlying distribution of data.

- Likelihood

$$l(\theta | X) = p(X | \theta) = \prod_{i=1}^n p(x_i | \theta)$$

→ independent

- Log Likelihood

$$L(\theta | X) = \log l(\theta | X) = \sum_i \log p(x_i | \theta)$$

- Maximum Likelihood Estimator

$$\hat{\theta} = \underset{\theta}{\operatorname{arg\,max}} L(\theta | X)$$

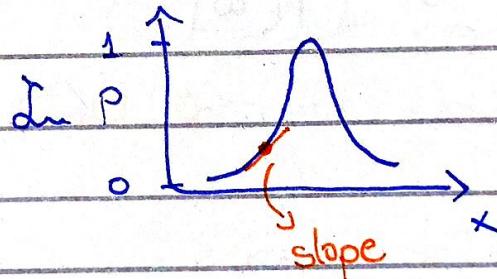
Machine Learning

- Bernoulli ; $P(x) = p^x (1-p)^{1-x}$

Log Likelihood $\{ \ln(p(x)) = \sum_t \log p^{x_t} (1-p)^{1-x_t} \}$ Important
(Probably wrong)

$$\text{Mean} = E(x) = \sum_x x p(x) = 1 \cdot p + 0 \cdot (1-p) \\ = \underline{p}$$

$$\text{Var} = \text{VAR}(x) = \sum_x (x - E(x))^2 p(x) = \underline{p(1-p)}$$



$$\frac{\partial \mathcal{L}}{\partial p} = 0 \Rightarrow \text{MLE } \hat{p} = \underbrace{\sum_t x_t}_{N} \frac{1}{N}$$

→ For multinomial, $K \geq 2$ states ;

$$x_i^t = \begin{cases} 1 & \text{if exp t chooses state i} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{MLE} \Rightarrow \hat{p}_i = \frac{\sum_t x_i^t}{N} \xrightarrow{\text{indusion of state i}}$$

Machine Learning

Parametric Classification

$$P(C_i | x) = \frac{P(C_i) P(x|C_i)}{\sum_k P(C_k) P(x|C_k)}$$

$$g_i(x) = \log(P(c_i)) + \log(P(x|c_i))$$

- If distribution is Gaussian;

$$\left[g_i(x) = -\frac{1}{2} \log(2\pi) - \log \sigma_i \dots \right. \\ \left. - \frac{(x - \mu_i)^2}{2\sigma_i^2} + \log(P(c_i)) \right]$$

→ Given $X = \{x^t, r^t\}_{t=1}^N$

ML Estimates \Rightarrow "estimates"

$$\hat{P}(c_i) = \frac{\sum r_i^t}{N} \rightarrow m_i = \frac{\sum x^t r_i^t}{\sum r_i^t}$$

$$s_i^2 = \frac{\sum (x^t - m_i)^2 r_i^t}{\sum r_i^t} \quad \begin{matrix} \downarrow \\ \text{Variance estimate} \end{matrix} \quad \begin{matrix} \downarrow \\ \text{mean estimate} \end{matrix}$$

- Actual mean $\sim E(x)$, where x is an RV expectation

$g_i(x) \rightarrow$ excluding constants

$$\hookrightarrow \underline{-(x - m_i)^2}$$

Choose C_i if $[x - m_i] = \min_k [x - m_k]$

\rightarrow suppose $p(x|C_1)$ and $p(x|C_2)$ are gaussians; the best discriminant point is the inter-crossing of the two gaussians.

- When the means are different and the variances are identical, there will only be one discriminant point.
- But if variances are different as well, there will be two discriminant points.
- Choose the discriminant with higher area in a particular region.

Exercise (Class)

$$g_1(x) = -\frac{1}{2} \log(2\pi) - \log S_1 - \frac{(x - m_1)^2}{2S_1^2} + \log P(C_1)$$

$$g_2(x) = -\frac{1}{2} \log(2\pi) - \log S_2 - \frac{(x - m_2)^2}{2S_2^2} + \log P(C_2)$$

Let $S_1 = S_2 = S_0$

$$\Rightarrow g_1(x) = g_2(x)$$

$$-\log \vec{s}_0 - \frac{(x - m_1)^2}{2s_0^2} = -\log \vec{s}_0 - \frac{(x - m_2)^2}{2s_0^2}$$

$$-(x - m_1)^2 = -(x - m_2)^2$$

$$-(x^2 - 2xm_1 + m_1^2) = -(x^2 - 2xm_2 + m_2^2)$$

$$-x^2 + 2xm_1 - m_1^2 = -x^2 + 2xm_2 - m_2^2$$

$$2x(m_1 - m_2) = m_1^2 - m_2^2$$

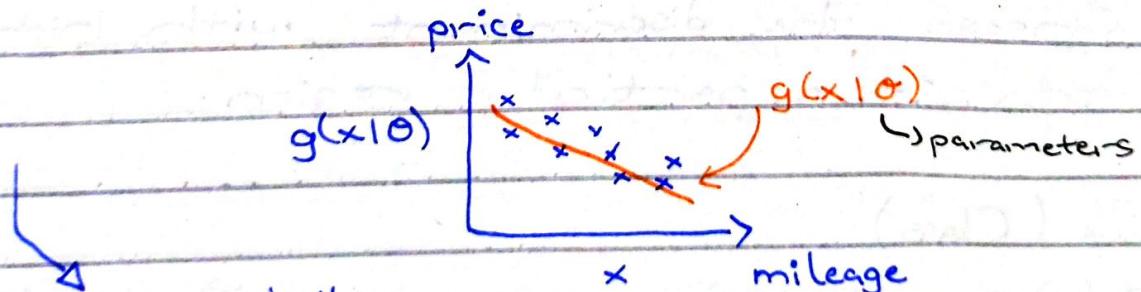
$$x = \frac{m_1^2 - m_2^2}{2(m_1 - m_2)}$$

$$x = \frac{1}{2}(m_1 + m_2)$$

Regression

- $r = f(x) + \epsilon$

$$\epsilon \sim N(0, \sigma^2)$$



now, we don't

know what the

class likelihood $P(x|C)$

is ;

mean is our estimator

$$p(r|x) \sim N(g(x|\theta), \sigma^2)$$

$$\mathcal{L}(\theta, x) = \log \prod_t P(x^t, r^t)$$

$$\therefore P(x, y) = P(y|x) P(x)$$

Linear Regression

$$g(x^t | w_1, w_0) = w_1 x^t + w_0$$

$$\sum_t r^t \quad \text{sum of labels} = Nw_0 + w_1 \sum_t x^t \quad (\text{i})$$

$$x^t = \begin{bmatrix} 1 \\ x^t \\ \vdots \\ N \end{bmatrix}$$

$$\sum_t r^t x^t = w_0 \sum_t x^t + w_1 \sum_t (x^t)^2 \quad (\text{ii})$$

$$A = \begin{bmatrix} N & \sum_t x^t \\ \sum_t x^t & \sum_t (x^t)^2 \end{bmatrix} \quad w = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \quad y = \begin{bmatrix} \sum_t r^t \\ \sum_t r^t x^t \end{bmatrix}$$

$$w = A^{-1}y$$

non-iterative,
mathematical solution

17/10/23

Machine Learning

Linear Regression

$$X = [x^t]_{t=1}^N, \quad W = [w_0 \ w_1]$$

$$g(x^t | w_1, w_0)_{t=1}^N = (w_1 x^t + w_0)_{t=1}^N$$

Polynomial Regression

$$g(x^t | w_k, w_{k-1}, \dots, w_1, w_0) = w_k x^t + \dots + w_1 x^t + w_0$$

~ After solving :

$$\hat{w} = (D^T D)^{-1} D^T r$$

Multivariate Data

$$X = \begin{bmatrix} x_1^1 & \dots & x_d^1 \\ x_1^2 & \dots & x_d^2 \\ \vdots & & \vdots \\ x_1^N & \dots & x_d^N \end{bmatrix} \quad \text{Row } = x^1$$

$$\mu = [m_1 \ \dots \ m_d]^T \quad \text{Column wise mean}$$

$$\text{Covariance : } \sigma_{ij} = \text{Cov}(x_i, x_j)$$

$X - \mu$ { Replicate μ such that row-wise subtraction occurs }

$$\text{Cov} = E[(x - \mu)^T (x - \mu)]$$

$$\text{Corr} = \rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$$

- If Missing Values exist, take column wise mean and replace the value missing
 - [can also take mode]
 - [can also ignore {not suitable if dataset size is small}]

$$x \sim N(\mu, \Sigma)$$

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

→ shapes : $\underbrace{x (d \times 1), \mu (d \times 1), \Sigma (d \times d)}$
 Result as above
 is (1×1)

$$\gg p(x^1) = \text{scalar} \\ \vdots \\ p(x^N) = \text{scalar} \quad \left. \right\} N \text{ values}$$

Machine Learning

- When features are independent along with samples:
 - auto covariance (diagonal) remain non-zero
 - off-diagonals (cross covariance) are zero
- (Naive Bayes)

$$p(x) = \prod_{i=1}^d p_i(x_i)$$

↳ (Joint probability when independence)

Dimensionality Reduction

Curse of Dimensionality { when $d > N$ }

$$\begin{matrix} & | & \dots & d \\ \vdots & [& x_1' & \dots & x_d'] \\ N & [& x_1^N & \dots & x_d^N] \end{matrix}$$

N should be at least 3 times more than d for there to be reasonable decisions

Principal Component Analysis (PCA)

↳ simple yet effective

Machine Learning

$y = mx + c$; $m = r \frac{\sigma_y}{\sigma_x}$ where r = correlation coefficient

$$r = \frac{[\hat{x} - \bar{x}][\hat{y} - \bar{y}]}{\sqrt{(\hat{x} - \bar{x})^2(\hat{y} - \bar{y})^2}}$$

$$\hat{y} = m\hat{x} + c; c = \bar{y} - m\bar{x}$$

$$\hat{y} = m(\hat{x} - \bar{x}) + \bar{y}$$

Principal Component Analysis

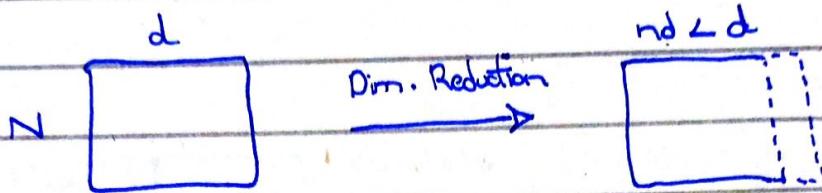
- └ Find mean of features
- └ Subtract for mean centering
- └ PC1 [Principal Component] passes through mean in -towards the largest variance
- └ PC2 is orthogonal to PC1; Project the 3D points onto the 2D plane

PCA gives unique { will result in
 (independent) features { features that do
 not intermingle

Variance

- The largest value in the data is represented by largest eigenvalue.

Machine Learning



In PCA, covariance is for identifying those unique eigenvalue to yield highest class deentangling.
Dimensionality Reduction

Linear Discriminant Analysis (LDA)

↳ makes use of class information

Steps:

↳ Find the separation of the classes

$$S_B = \sum_{i=1}^g N_i (\bar{x}_i - \bar{\bar{x}})(\bar{x}_i - \bar{\bar{x}})^T$$

↙ between class scatter matrix

↳ Find the separation within the classes

$$S_W = \sum_{i=1}^g (N_i - 1) S_i$$

$$= \sum_{i=1}^g \sum_{j=1}^{N_i} (x_{ij} - \bar{x}_i)(x_{ij} - \bar{x}_i)^T$$

L Maximize b/w class variance

Minimize in-class variance

Neighbourhood Component Analysis

[Given $\mathbf{X} [N \times d]$]

[Initialize mapping matrix $\mathbf{A} [d \times nd]$]

• where $nd = d$ new dimension

L on transformed space, find stochastic
nearest neighbours

$$p_{ij} = \frac{\exp(-\|\mathbf{A}\mathbf{x}_i - \mathbf{A}\mathbf{x}_j\|^2)}{\sum_{k \neq i} \exp(-\|\mathbf{A}\mathbf{x}_i - \mathbf{A}\mathbf{x}_k\|^2)}, \quad \underline{p_{ii} = 0}$$

• Probability that point j, the
neighbour of point i belongs to
the same class as i.

$$\underline{p_i = \sum_{j \in C_i} p_{ij}}$$

$$\underline{f(\mathbf{A}) = \sum_i p_i} \quad \rightarrow \frac{\partial f}{\partial \mathbf{A}} = 0$$

\rightarrow & find optimum \mathbf{A}

14/11/23

Machine Learning

Gaussian Mixture Method {GMM} Model

$d \triangleq$ dimensionality

$K \triangleq$ number of Gaussians

→ GMM is a weighted set of Gaussians summed to form a single distribution.

$$p(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k),$$

with $\pi_k \triangleq$ the mixing coefficients.

$$\boxed{\sum_{k=1}^K \pi_k = 1}$$

Constraint

and $\pi_k > 0 \forall k$

of the form $p(x|\theta)$

$$\left\{ \begin{array}{l} \text{log-likelihood ; } L = \ln p(x | \pi, \mu, \Sigma) \\ = \sum_n \ln \left(\sum_k \pi_k N(x^{(n)}; \mu_k, \Sigma_k) \right) \\ \hookrightarrow n \text{ as } X \text{ is a vector of all data points} \\ \text{wrt. } \Theta = \{\pi_k, \mu_k, \Sigma_k\} \end{array} \right.$$

$x^{(n)} \in X$

• Find Θ using:

$$\frac{\partial L}{\partial \pi_k} = 0 ; \frac{\partial L}{\partial \mu_k} = 0 ; \frac{\partial L}{\partial \Sigma_k} = 0$$

~ in-class

$$p(x) = \sum_k \pi_k N(x | \mu_k, \Sigma_k)$$

Let z_k be the clusters ; \sim gaussians

$$p(x) = \sum_k p(x, z=k) = \sum_k \underbrace{p(x|z=k)}_{N(x|\mu_k, \Sigma_k)} \underbrace{p(z=k)}_{\pi_k}$$

Take $\ln \rightarrow$ Solve and we will
get expression of L

Machine Learning

Linear Regression

$$E = \frac{1}{2} \sum_n [r^{(n)} - w_0 x^n - w_1]^2$$

$$\frac{\partial E}{\partial w_0} = \frac{2}{2} \left[\sum_n (r^{(n)} - w_0 x^n - w_1) \right] [-w_0 N] = 0$$

$\underbrace{\sum_{t=1}^n r^t}_{\text{but}} = N w_0 + w_1 \sum_t x^{(t)}$ — i

$$\frac{\partial E}{\partial w_1} = \frac{2}{2} \left[\sum_t (r^{(t)} - w_0 x^{(t)} - w_1) \right] \left[-\sum_t x^{(t)} \right] = 0$$

$$\sum_t r^{(t)} x^{(t)} = w_0 \sum_t x^{(t)} + w_1 \sum_t [x^{(t)}]^2 \quad \text{— ii}$$

Multivariate Case

Model $d = \text{features}$

$$g(x | w_0, w_1, \dots, w_d) = w_0 + w_1 x_1 + \dots + w_d x^d$$

The rest is the same as $w = (X^T X)^{-1} X^T r$

Moore-Penrose
Inverse

CNMM $p(x) = \sum_{k=1}^K \pi_k N(x | \mu_k, \Sigma_k)$

Maximum height of $p(x)$ is one.

$$\therefore \sum_{k=1}^K \pi_k = 1$$

$$\hookrightarrow p(x) = \prod_n \left(\sum_k \pi_k N(x^{(n)} | \mu_k, \Sigma_k) \right)$$

^{log} likelihood $L = \sum_n \ln \left(\sum_k \pi_k N(x^{(n)} | \mu_k, \Sigma_k) \right)$

$$\frac{\partial L}{\partial \pi_k} = 0 ; \frac{\partial L}{\partial \mu_k} = 0 ; \frac{\partial L}{\partial \Sigma_k} = 0$$

e f g

» 3k equations

» Latent Variable

evidence

$$p(x) = \sum_{k=1}^K p(x, z=k) = \sum_{k=1}^K \underbrace{\pi_k}_{\text{evidence}} \underbrace{p(z=k) p(x|z=k)}_{N(x|\mu_k, \Sigma_k)}$$

» Parameters

$$\mu_k = \frac{\sum_{m=1}^M \mathbb{1}_{[z^{(m)}=k]} x^{(m)}}{M} = \frac{\sum_{m=1}^M \frac{1}{M} \mathbb{1}_{[z^{(m)}=k]} x^{(m)}}{\sum_{m=1}^M \mathbb{1}_{[z^{(m)}=k]}}$$

$$\pi_k = \frac{1}{N} \sum_{m=1}^N \mathbb{1}_{[z^{(m)}=k]}$$

$$\Sigma_k = \frac{\sum_{m=1}^N \mathbb{1}_{[z^{(m)}=k]} (x^{(m)} - \mu_k)(x^{(m)} - \mu_k)^T}{\sum_{m=1}^N \mathbb{1}_{[z^{(m)}=k]}}$$

Expectation - Maximization (EM)

(EM) Algorithm

$$P(C|x) = \frac{P(c)p(x|c)}{P(x)}$$

E Step

Compute posterior probability of / over z
with respect to k } responsibility
initial model } fancy way to say
posterior $\hat{\pi}_k$

M Step

Maximize the probability
that it would generate the data it
is responsible for.

$$\begin{aligned} \text{E step } \left\{ \hat{\pi}_k &= p(z=k|x) = \frac{p(x|z=k) p(z=k)}{p(x)} \right. \\ &= \left. \frac{\pi_k N(x|\mu_k, \Sigma_k)}{\sum_j^K \pi_j N(x|\mu_j, \Sigma_j)} \right. \end{aligned}$$

$$\text{M step } \left\{ \mu_k, \pi_k, \Sigma_k \right.$$

→ Equations from slides

use M-step parameters for evaluating
Log-likelihood for convergence

28/11/23

Machine Learning

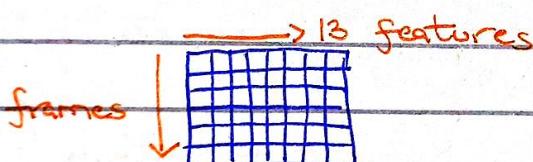
Speaker Recognition

10 files 10 files ... 10 files
[S₁] [S₂] ... [S₁₀]

[For Each Speaker →

10 files → MFCC .npy → 0.2 train / test split ↗ cont.
• cont. split [train (8 files)
 test (2 files)

MFCC { Mel frequency cepstral coefficients }



Perceptron

→ like num, just do slides
did the same shit in CV

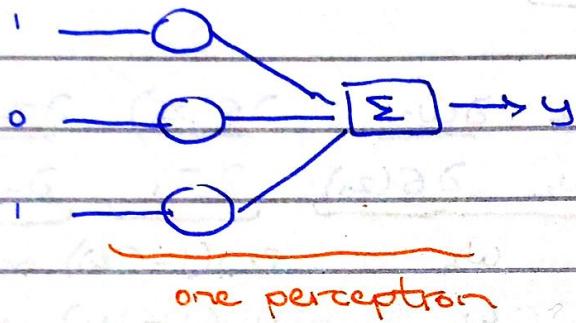
bias := offset parameter → weights can not
trainable control offsets
parameter as well

4/12/23

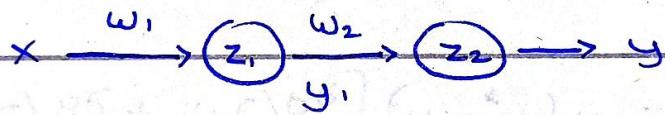
Machine Learning

$$w_{\text{new}} := w_{\text{old}} + \alpha (\text{desired} - \text{output}) / \text{input}$$

\hookrightarrow learning rate



multi-layer perceptron [one hidden layer]



$$z_1 = xw_1$$

$$y_1 = f(z_1) = \sigma(z_1) \quad f(\cdot) = \text{sigmoid}$$

$$z_2 = w_1 y_1 w_2$$

$$y = f(z_2) = \sigma(z_2) = \sigma(\sigma(xw_1)w_2)$$

$$\text{Let } \epsilon \text{ error (MSE)} = \frac{1}{2} [y^* - y]^2$$

\hookrightarrow desired / target

$$\frac{\partial \epsilon}{\partial w_1} = ? \quad \frac{\partial \epsilon}{\partial w_2} = ?$$

$$\gg \frac{\partial \epsilon}{\partial w_j} = - (y^* - y) \frac{\partial y}{\partial w_j} - A$$

$$\begin{aligned}
 \frac{\partial y}{\partial w_2} &= \frac{\partial \sigma(z_2)}{\partial z_2} \frac{\partial z_2}{\partial w_2} \\
 &= \frac{\partial \sigma(z_2)}{\partial z_2} \frac{\partial(y, w_2)}{\partial w_2} = \sigma(z_2)[1 - \sigma(z_2)] y, \\
 &= \underbrace{y(1-y)}_{\delta_2} y_1 - B
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial y}{\partial w_1} &= \frac{\partial \sigma(z_2)}{\partial z_2} \frac{\partial z_2}{\partial w_2} \frac{\partial w_2}{\partial \sigma(z_1)} \frac{\partial \sigma(z_1)}{\partial z_1} \frac{\partial z_1}{\partial w_1}, \\
 &\quad \sigma(z_2)(1-\sigma(z_2)) \quad w_2 \quad \sigma(z_1)(1-\sigma(z_1)) \times \\
 &\quad y(1-y) \quad w_2 \quad y_1(1-y_1) \times \\
 &= \underbrace{y(1-y) w_2 y_1(1-y_1)}_{\delta_1} \times - C
 \end{aligned}$$

$$-A : \frac{\partial e}{\partial w_j} = -(y^* - y) \left[\frac{\partial y}{\partial w_1} + \frac{\partial y}{\partial w_2} \right]$$

$$\circ w_2 := w_2 - \alpha \delta_2 y_1$$

$$\circ w_1 := w_1 - \alpha \delta_1 \times$$

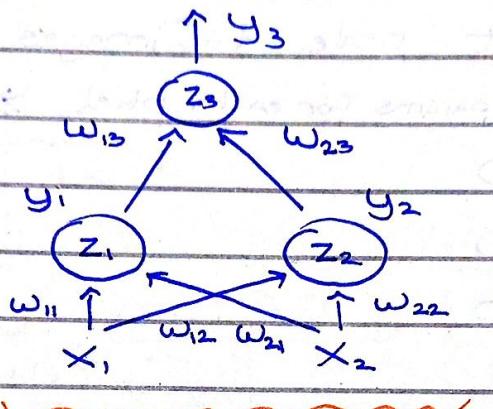
12/11/23

Machine Learning

Training a NN

Key

desired : y^*
activation
 sigmoid



$$\begin{cases} y_1 = \sigma(z_1) \\ z_1 = x_1 w_{11} + x_2 w_{21} \\ y_2 = \sigma(z_2) \\ z_2 = x_1 w_{12} + x_2 w_{22} \\ y_3 = \sigma(z_3) \\ z_3 = y_1 w_{13} + y_2 w_{23} \end{cases}$$

$$E = \frac{1}{2} [y^* - y_3]^2 \quad \left[\begin{array}{l} n \text{ neurons will generate} \\ n \text{ Ss} \end{array} \right]$$

$$\frac{\partial E}{\partial w} = ?$$

$$y \equiv y_3$$

Output Layer $\frac{\partial E}{\partial w_{i3}} = \frac{\partial E}{\partial w_{13}} + \frac{\partial E}{\partial w_{23}}, i = 1, 2$

$$\frac{\partial E}{\partial w_{i3}} = -(y^* - y) \left[\frac{\partial y}{\partial w_{13}} + \frac{\partial y}{\partial w_{23}} \right]$$

$$\frac{\partial y}{\partial z_3} \frac{\partial z_3}{\partial w_{i3}}$$

$$y_i(1-y_i) y_i$$

δ_3

$$\frac{\partial y}{\partial z_3} \frac{\partial z_3}{\partial w_{i3}}$$

$$y_3(1-y_3) y_2$$

$$\frac{\partial E}{\partial w_{i3}} = -(y^* - y) y_3 (1-y_3) [y_1 + y_2]$$

Hidden Layer

$$\frac{\partial E}{\partial w_{ij}} = -(y^* - y) \left[\underbrace{\frac{\partial y}{\partial w_{12}}}_{\frac{\partial y_3}{\partial z_3} \frac{\partial z_3}{\partial y_2} \frac{\partial y_2}{\partial z_2} \frac{\partial z_2}{\partial w_{12}}} + \underbrace{\frac{\partial y}{\partial w_{11}}}_{\frac{\partial y_3}{\partial z_3} \frac{\partial z_3}{\partial y_1} \frac{\partial y_1}{\partial z_1} \frac{\partial z_1}{\partial w_{11}}} \right]$$

$$\frac{\partial y_3}{\partial z_3} \frac{\partial z_3}{\partial y_2} \frac{\partial y_2}{\partial z_2} \frac{\partial z_2}{\partial w_{12}} \quad \frac{\partial y_3}{\partial z_3} \frac{\partial z_3}{\partial y_1} \frac{\partial y_1}{\partial z_1} \frac{\partial z_1}{\partial w_{11}}$$

$$y_3(1-y_3) w_{23} y_2(1-y_2) x_1 \quad y_3(1-y_3) w_{13} y_1(1-y_1) x_1$$

$$\frac{\partial E}{\partial w_{ij}} = -(y^* - y) \sim \text{OKE Ignore}$$

Sir's way \approx \hat{s}_j $\hat{o}_j =$

$$\frac{\partial E}{\partial w_{12}} = -(y^* - y) \left[\underbrace{\frac{\partial y}{\partial w_{12}}}_{\frac{\partial y_3}{\partial z_3} \frac{\partial z_3}{\partial y_2} \frac{\partial y_2}{\partial z_2} \frac{\partial z_2}{\partial w_{12}}} + \underbrace{\frac{\partial y}{\partial w_{22}}}_{\frac{\partial y_3}{\partial z_3} \frac{\partial z_3}{\partial y_2} \frac{\partial y_2}{\partial z_2} \frac{\partial z_2}{\partial w_{22}}} \right]$$

$$\frac{\partial y_3}{\partial z_3} \frac{\partial z_3}{\partial y_2} \frac{\partial y_2}{\partial z_2} \frac{\partial z_2}{\partial w_{12}} \quad \frac{\partial y_3}{\partial z_3} \frac{\partial z_3}{\partial y_2} \frac{\partial y_2}{\partial z_2} \frac{\partial z_2}{\partial w_{22}}$$

$$y_3(1-y_3) w_{23} y_2(1-y_2) x_1 \quad y_3(1-y_3) w_{23} y_2(1-y_2) x_2$$

$$\frac{\partial E}{\partial w_{12}} = \frac{-(y^* - y) y_3(1-y_3) w_{23} y_2(1-y_2)}{\delta_2} [x_1 + x_2]$$

$$\text{Similarly, } \frac{\partial E}{\partial w_{11}} = \frac{-(y^* - y) y_3(1-y_3) w_{13} y_1(1-y_1)}{\delta_1} [x_1 + x_2]$$

Updates

$$\begin{aligned} w_{13} &= w_{13} - \gamma \delta_3 y_1 \\ w_{23} &= w_{23} - \gamma \delta_3 y_2 \\ &\vdots \\ w_n &= w_n - \gamma \delta_n y_n \end{aligned} \quad \left. \begin{array}{l} \text{The rest is same as} \\ \text{well} \end{array} \right\}$$

new old lr error inputs

Weight update equation decides the trajectory of E-curves. { should avoid local minima }

Linear activation := Identity (\circ) func.

$$x_{in} \xrightarrow{\theta} z \rightarrow y \quad y = z$$

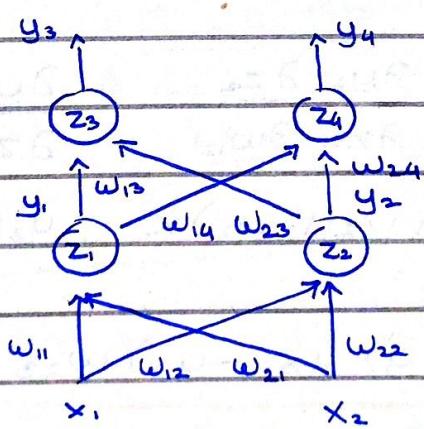
$z = \theta x_{in}$

More Practice

$$E = \frac{1}{2} [(y_3^* - y_3)^2 + (y_4^* - y_4)^2]$$

Activation:

Sigmoid



Weight Updates

- $w_{11} := w_{11} - \gamma \delta_1 x_1$ • $w_{12} := w_{12} - \gamma \delta_2 x_1$
- $w_{22} := w_{22} - \gamma \delta_2 x_2$ • $w_{21} := w_{21} - \gamma \delta_1 x_2$
- $w_{13} := w_{13} - \gamma \delta_3 y_1$ • $w_{14} := w_{14} - \gamma \delta_4 y_1$
- $w_{23} := w_{23} - \gamma \delta_3 y_2$ • $w_{24} := w_{24} - \gamma \delta_4 y_2$

not sure about
the negative sign

$$\begin{aligned} S_4 &= -y_4(1-y_4)(y_4^* - y_4) \\ S_3 &= -y_3(1-y_3)(y_3^* - y_3) \\ S_2 &= y_2(1-y_2)S_3 w_{23} \\ &\quad y_2(1-y_2)S_4 w_{24} \\ &= y_2(1-y_2)[S_3 w_{23} + S_4 w_{24}] \end{aligned}$$

$$S_1 = y_1(1-y_1)[S_3 w_{13} + S_4 w_{14}]$$

S_1 and S_2 shouldn't have -ive but sir says otherwise
check directly

12/12/23

Machine Learning

→ Replace dot product of FNN with *

convolution

$$z_1 = x * w_1$$

$$y_1 = \sigma(z_1)$$

$$z_2 = y_1 * w_2$$

$$y_2 = \sigma(z_2)$$

$$z_2 = y_1 * w$$

why did i write it
again?

am i high?
or am i just drunk
in love? ^^\n

$$\rightarrow z_1 = \sum_{k=-\infty}^{\infty} x(k) w_1(k-n)$$

lmao no, i'm
emotionally dead

1/11/24

Machine Learning

On CNNs & Shapes

Input Layer : $5 \times 5 \times 3$ | width x height x chan.



Conv 2D Layer :
$$\frac{I - F + 2P}{S} + 1$$

P: padding

F: filter

I: image

/ Example
$$\begin{cases} I = 5 \times 5 \times 3 \\ F = 3 \times 3 \\ P = 2 \times 2 \times 2 \times 2 \\ S = 2 \end{cases}$$

1 $O = \frac{5 - 3 + 2(2)}{2} + 1 = 4$



Pooling Layer :
$$\frac{I - F}{S} + 1$$

→ only stride controls output
shape

Example
$$\begin{cases} F = 2 \times 2 \\ S = 2 \end{cases} \rightarrow \text{half}$$
 $S = 1 \rightarrow \text{no effect}$

Perform ceil() on decimal values

Global Average Pooling

↳ if Input: $7 \times 7 \times 3$ | $\xrightarrow{\text{average of all}}$ 3 is channels
Output: $1 \times 1 \times 3$ | $1 \times 1 \times \# \text{chan.}$

Batch Normalization

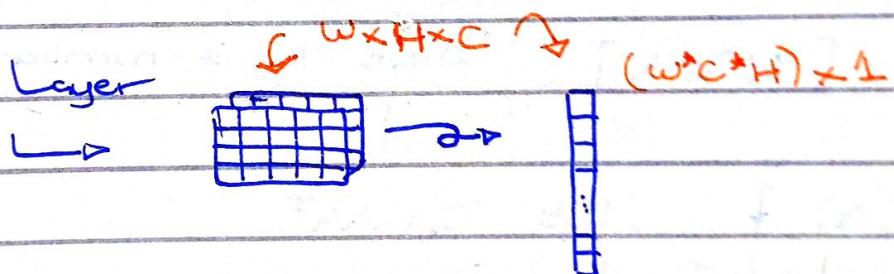
↳ No change in dimension

$$\text{Data}' = \frac{[\text{Data} - \mu_0]}{\sigma_0} \quad \begin{array}{|l} \text{Normalized} \\ \text{o mean} \\ \text{+ variance / STD} \end{array}$$

Dropout Layer

↳ random dropping

Flatten Layer



epochs $\propto 1 / \text{l.r.}$