

Machine Learning 1

Lecture 1 - Overview and Probability Theory

-Dr. Daniel Worrall-



*Slides created by:
Rianne van den Berg*

Image credit: Kirillm | Getty Images

Admin: grading policy

- ▶ 20 % homework assignments
- ▶ 20 % programming assignments
- ▶ 60 % final exam
- ▶ Passing requires at least **50%** of the points in the final exam!!!
- ▶ No midterm exam
- ▶ Resit exam replaces 100% of the grade

Admin: lab assignments

- ▶ Laptop exercise classes (laptop colleges) on Wednesdays
- ▶ 3 assignments to hand in
- ▶ Assignments available on canvas
- ▶ must be handed in in groups of 2.
- ▶ Late submissions will not be graded.
- ▶ For detailed instructions, see canvas!

Admin: homework assignments

- ▶ Exercise classes on Thursdays
- ▶ During class: study practice exercises
- ▶ At home: make hand-in assignments.
- ▶ Must be handed in individually.
- ▶ Must be written **in Latex**. Handwritten papers or scans are not allowed.
- ▶ Solutions should contain derivations.
- ▶ You can help each other and discuss, but do not copy-paste!

Admin: lectures

- ▶ Tuesdays and Wednesdays 17-19
- ▶ Keep an eye on datanose for irregularities in the schedule
- ▶ Slides will be made available on Canvas

Admin: general

p.d.force@uva.nl

- ▶ Please stick to your assigned group for laptop and exercise classes
- ▶ Questions: email your TA!
- ▶ If TA cannot answer your question: email me.
- ▶ See Canvas for email addresses

Prerequisite knowledge

- ▶ Probability Theory (+ basics of Statistics)

Bayesian Statistics for Machine Learning

- ▶ Programming – Python

- ▶ Linear Algebra

- ▶ Calculus

- ▶ Vector Calculus (+ some Optimisation)

Prerequisite knowledge

Extra resources:

- ▶ Self-study course Mathematics for Artificial intelligence
(see canvas)
- ▶ Book: Mathematics for Machine Learning (see canvas)
Chapter 5 "

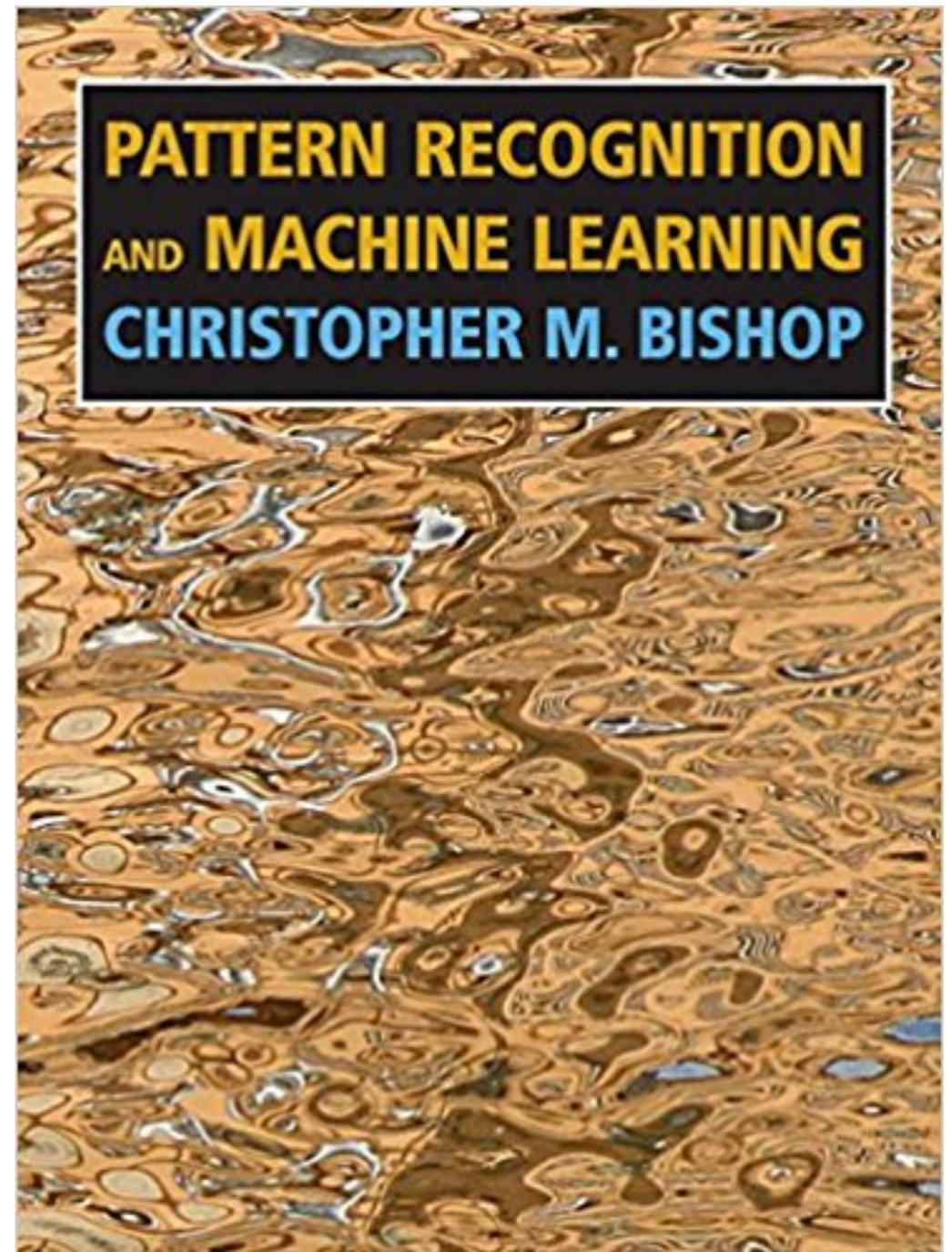
Literature: main book

Pattern Recognition and
Machine Learning

- Christopher Bishop

Machine learning from a
probabilistic point of view

Free online

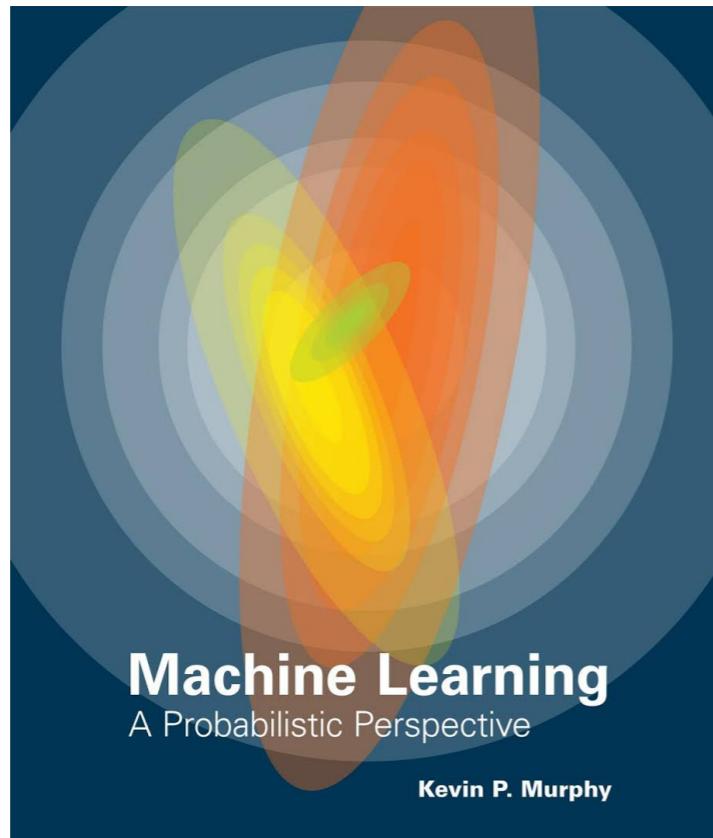


Literature: other books

Machine Learning - A Probabilistic Perspective

- Kevin P. Murphy

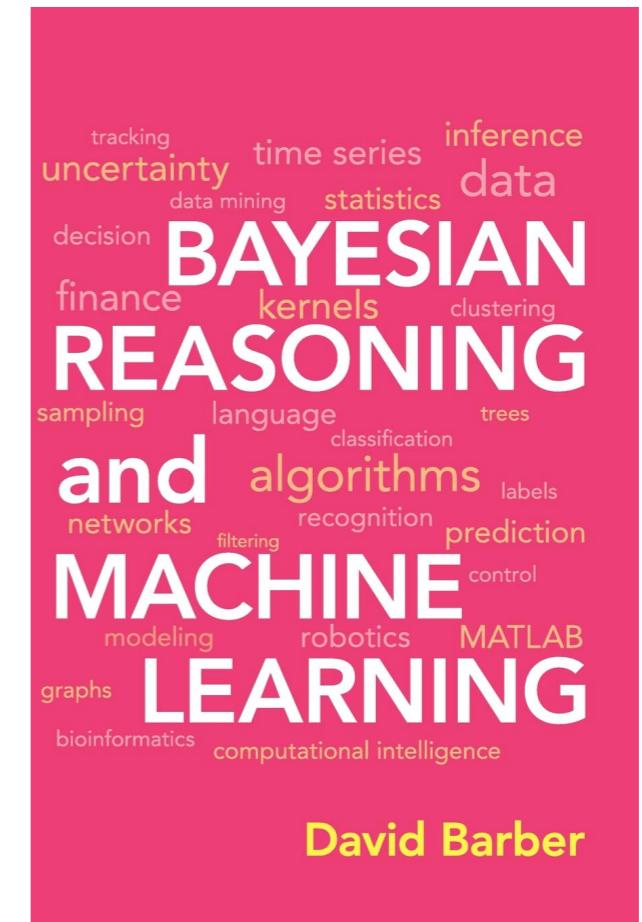
Best alternative to Bishop, maybe better



Bayesian Reasoning and Machine Learning

- David Barber

Machine learning viewed from Bayesian statistics

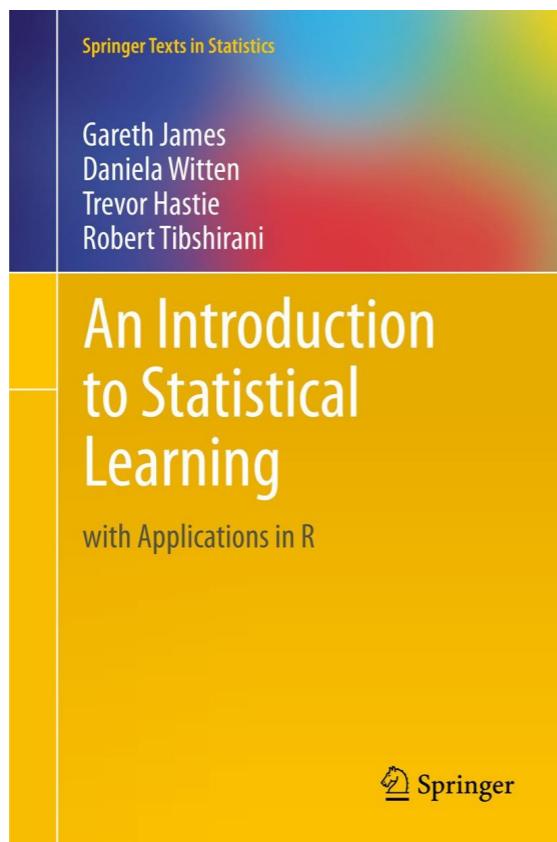


Literature: other books

An Introduction to Statistical Learning

- Gareth James, Daniela Witten, Trevor Hastie, Robert Tibshirani,

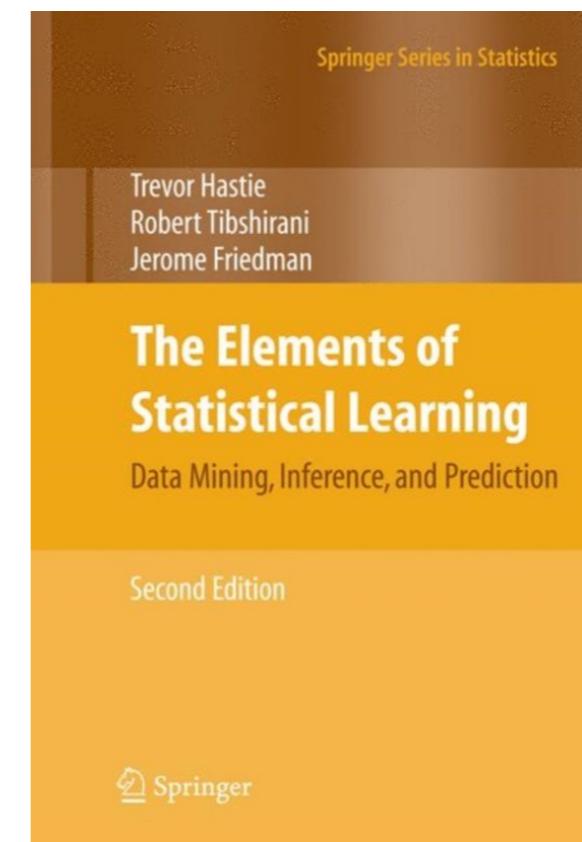
Introduction to Machine learning as a statistical tool.



The Elements of Statistical Learning

- Trevor Hastie, Robert Tibshirani, Jerome Friedman

More advanced view of Machine learning as a statistical tool.



Literature: Daniel's addition

Information Theory, Inference, and Learning Algorithms

- David MacKay

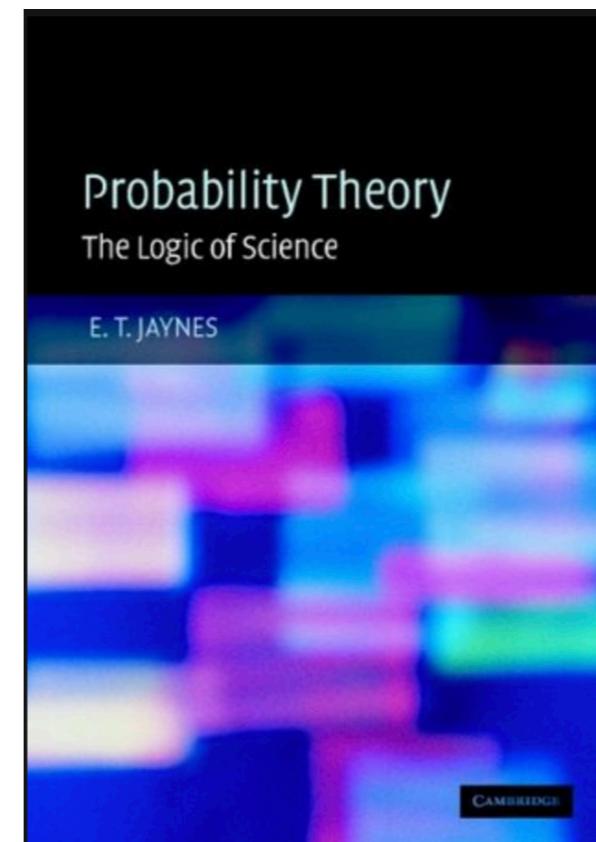
[https://www.inference.org.uk/itprnn/
book.pdf](https://www.inference.org.uk/itprnn/book.pdf)



Probability Theory: The Logic of Science

- E. T. Jaynes

A pure Bayesian treatment of probability theory



Questions?

Overview

1. Overview

1. What is machine learning?
2. Different types of machine learning
2. Recap: probability theory

What is Machine Learning?

Computer
powered

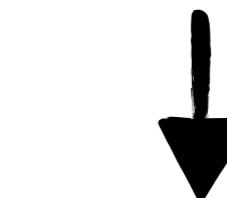


plus Engineering



Data driven AI
(in contrast to symbolic AI)

Improving on tasks
via processing data



aka Statistics



What is machine learning?

“A computer program is said to learn from experience **E** with respect to some class of tasks **T** and performance measure **P** if its performance at tasks in **T**, as measured by **P**, improves with experience **E**.”

*metric/
score*

- Tom M. Mitchell

data

Classification: Postcodes

28×28 pixels $\in \{0, \dots, 255\}$



2 = "2"
8 = "8"
0 = "0"

9 = "9" or "8"
8 = "8" or "9"

ground-truth
outputs (predictions)

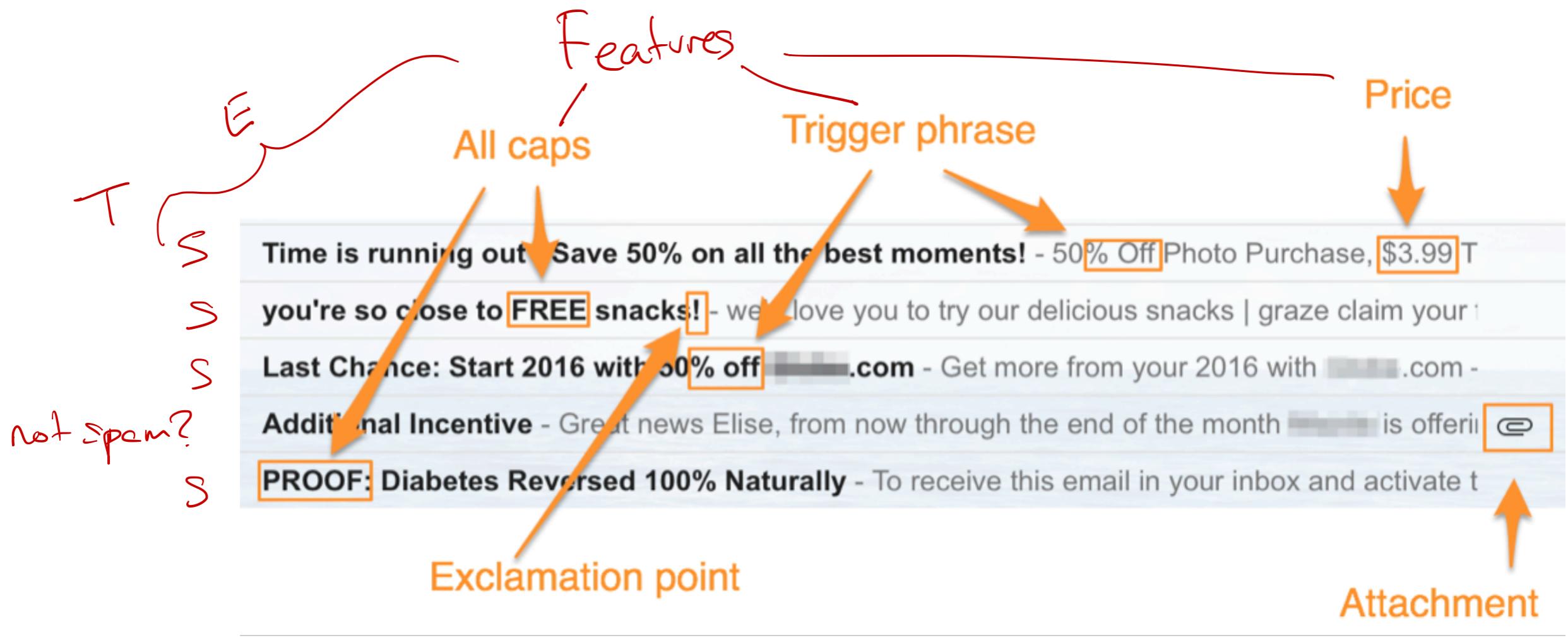
$$\text{accuracy}(y, \hat{y}) = \frac{1}{n_{\text{samples}}} \sum_{i=1}^{n_{\text{samples}}} \mathbb{1}[y_i = \hat{y}_i]$$

indicator func
 $= \begin{cases} 1 & \text{if } y_i = \hat{y}_i \\ 0 & \text{otherwise} \end{cases}$

label / aleatoric noise

MNIST dataset

Classification: Spam filtering



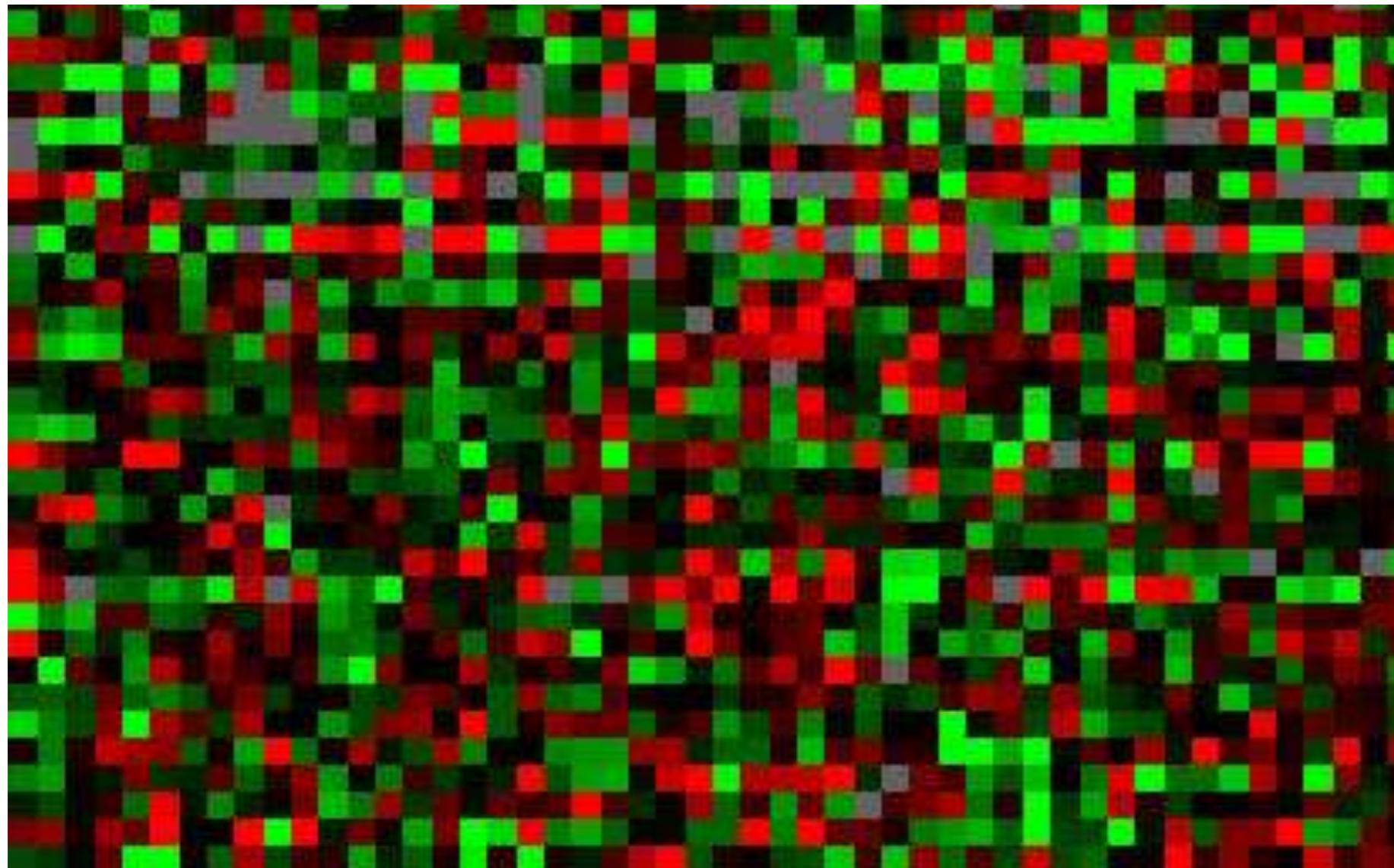
R
 $\text{accuracy}(y, \hat{y}) = \text{like on previous slide}$

Examples of spam emails. [source: Yesware]

Clustering: Genomics-Tumor data

E
J

☒ negative: underexpressed ☒ positive; overexpressed ☒ missing values

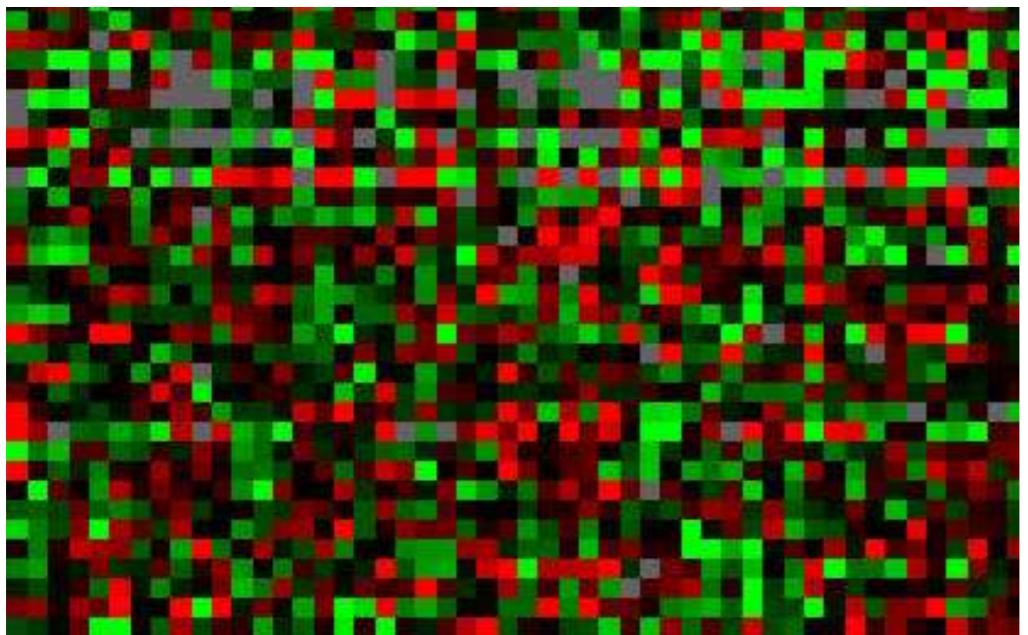


ESTs
SIDW486740
SMALLNUC
ESTs
SIDW366311
SIDW357197
SID52979
ESTs
SID43609
SIDW416621
ERLUMEN
TUPLE1TUP1
SIDW428642
SID381079
SIDW298052
SIDW417270
SIDW362471
ESTsChr.15
SIDW321925
SID380265
SIDW308182
SID381508
SID377133
SIDW365099
ESTsChr.10
SIDW325120
SID360097
SID375990
SIDW128368
SID301902
SID31984
SID42354

Gene

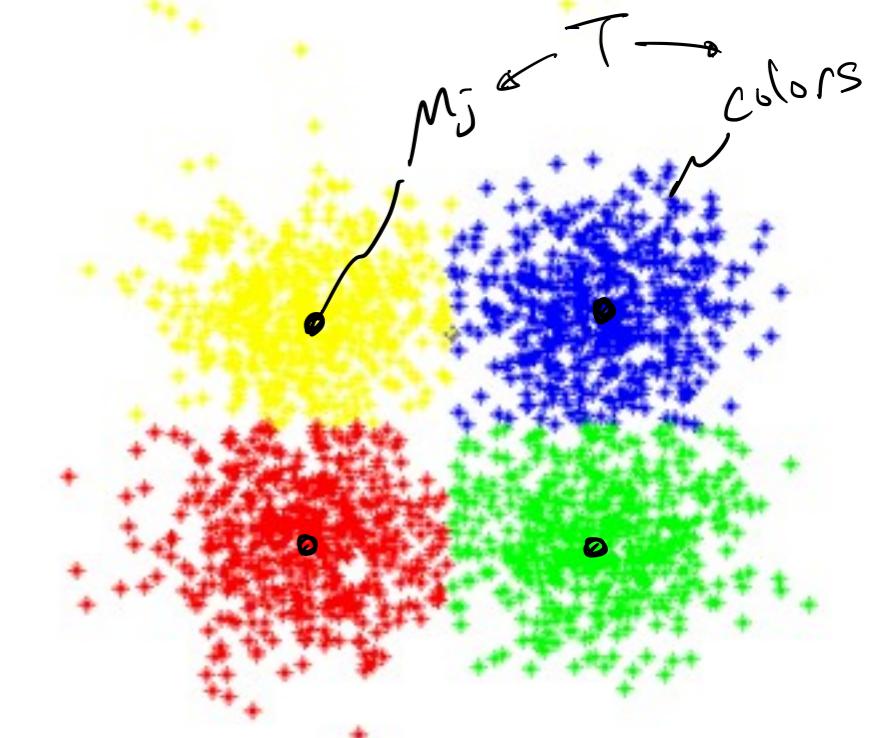
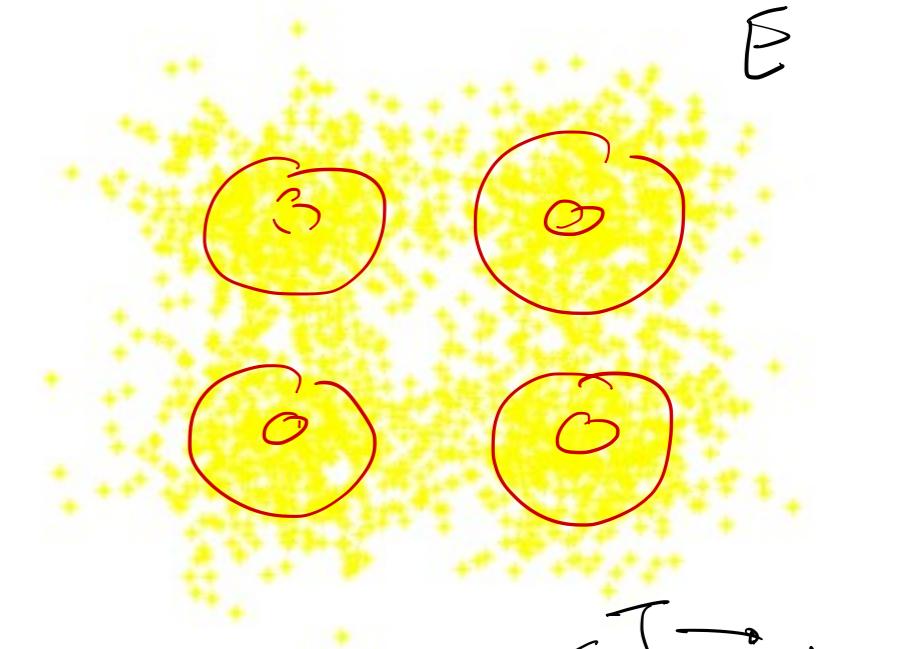
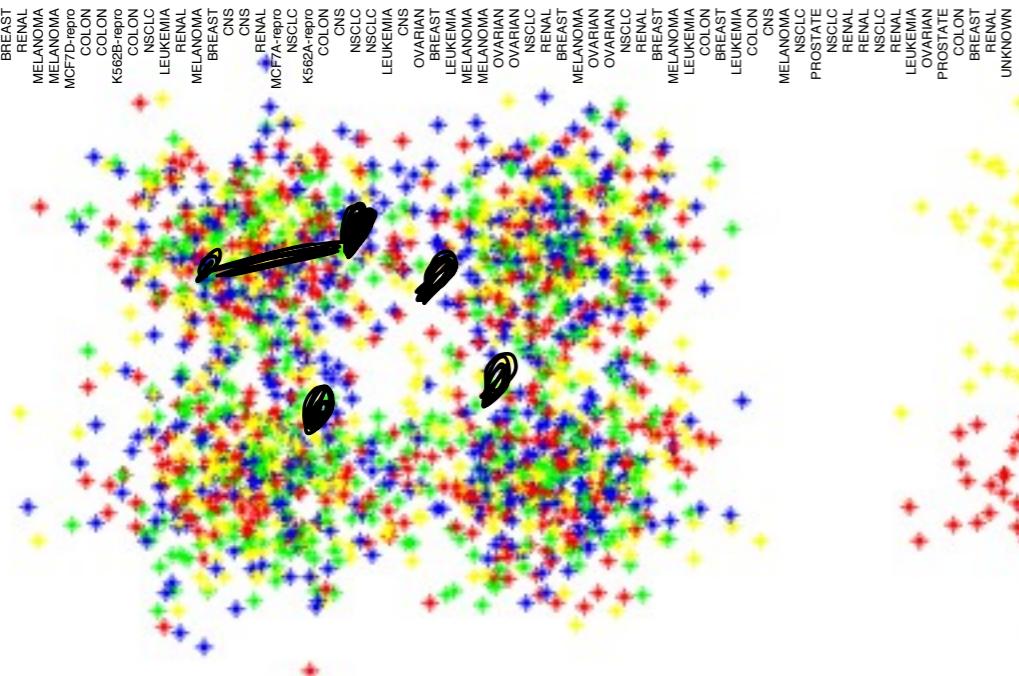
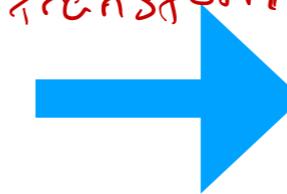
Expression matrix of genes (rows) for 64 human tumor samples (columns). [source: ESL 1.3]

Clustering: Genomics-Tumor data



ESTs
SIDW486740
SMALLNUC
ESTs
SIDW366311
SIDW357197
SID52979
ESTs
SID43609
SIDW416621
ERLUMEN
TUPLE1TUP-
SIDW428642
SID381079
SIDW298052
SIDW417270
SIDW362471
ESTsChr.15
SIDW321925
SID380265
SIDW308182
SID381508
SID377133
SIDW365099
ESTsChr.10
SIDW325120
SID360097
SID375990
SIDW128368
SID301902
SID31984
SID42354

feature transform



R

within cluster sum of squares =

$$\sum_{i=1}^{n_{\text{samples}}} \min_{\mu_j \in C} \|x_i - \mu_j\|_2^2$$

cluster means

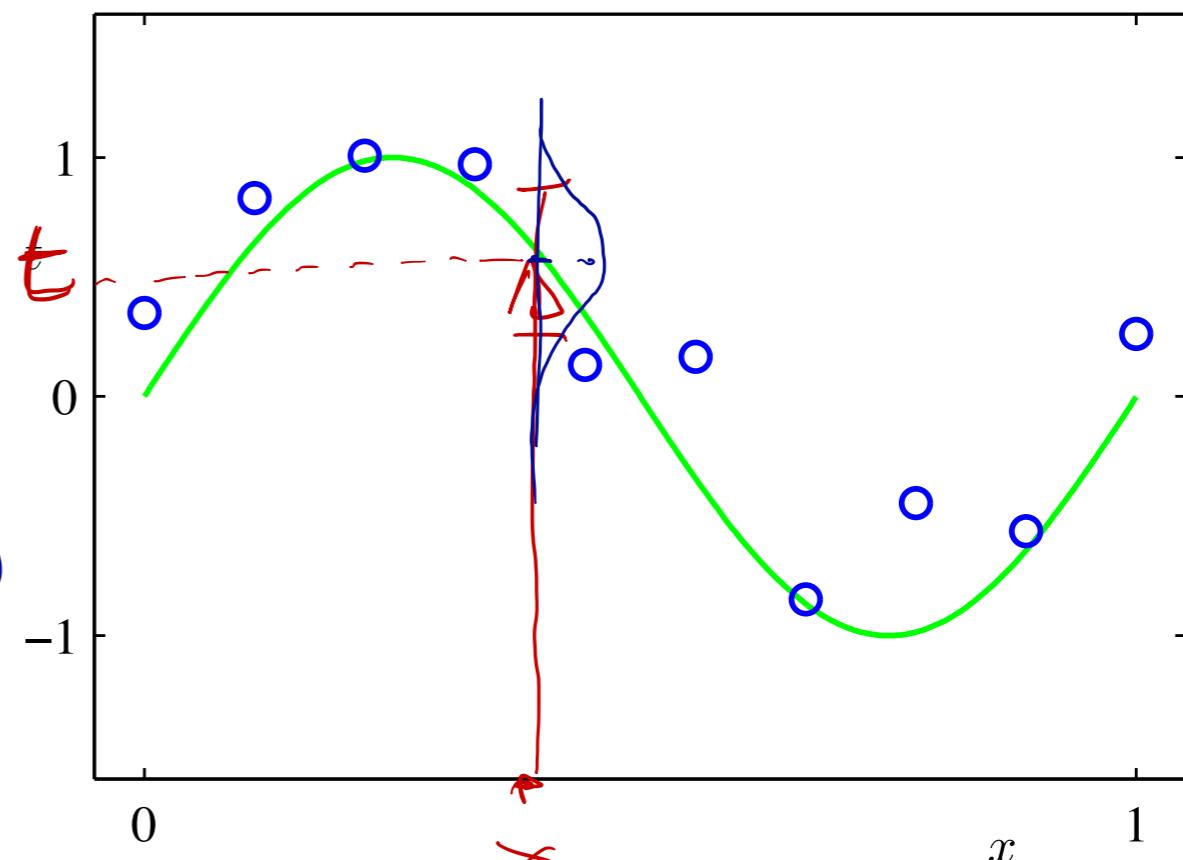
1D Regression

= curve fitting

input : x

target : $t = \sin(2\pi x) + \epsilon$

noise : $\epsilon \sim \mathcal{N}(\epsilon; 0, 1^2)$

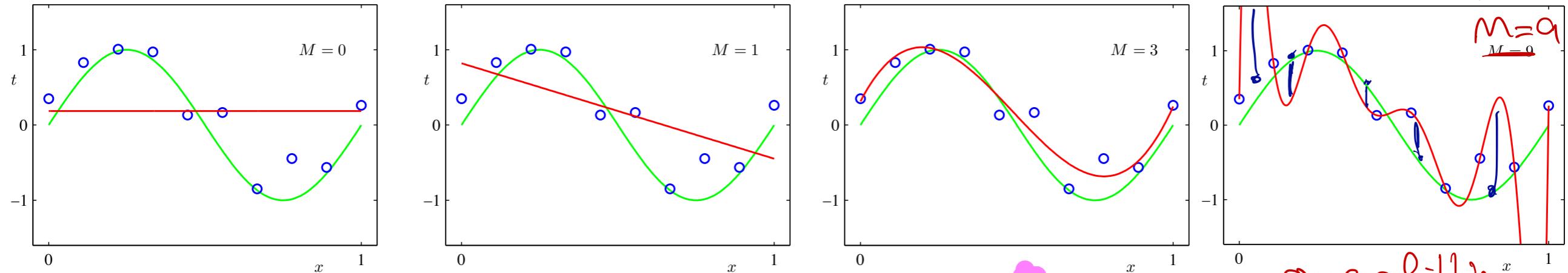


$p(t|x)$

1D Regression

$$f(x) = w_0 + w_1 x + w_2 x^2 + \dots + w_n x^n$$

$n = 9$



$$\text{MSE}(y, \hat{y}) = \frac{1}{n_{\text{samples}}} \sum_{i=1}^{n_{\text{samples}}} (y_i - \hat{y}_i)^2$$

Best performance on training set

Best performance on new datapoints

Q: On which datapoints should performance be measured?

Generalization: performance should be measured
on new data (test data)

Polynomials of order M (red) fit to data constructed as $t = \sin(2\pi x) + \epsilon$ (green)

Overview

1. Overview

1. What is machine learning?

2. Different types of machine learning

2. Recap: probability theory

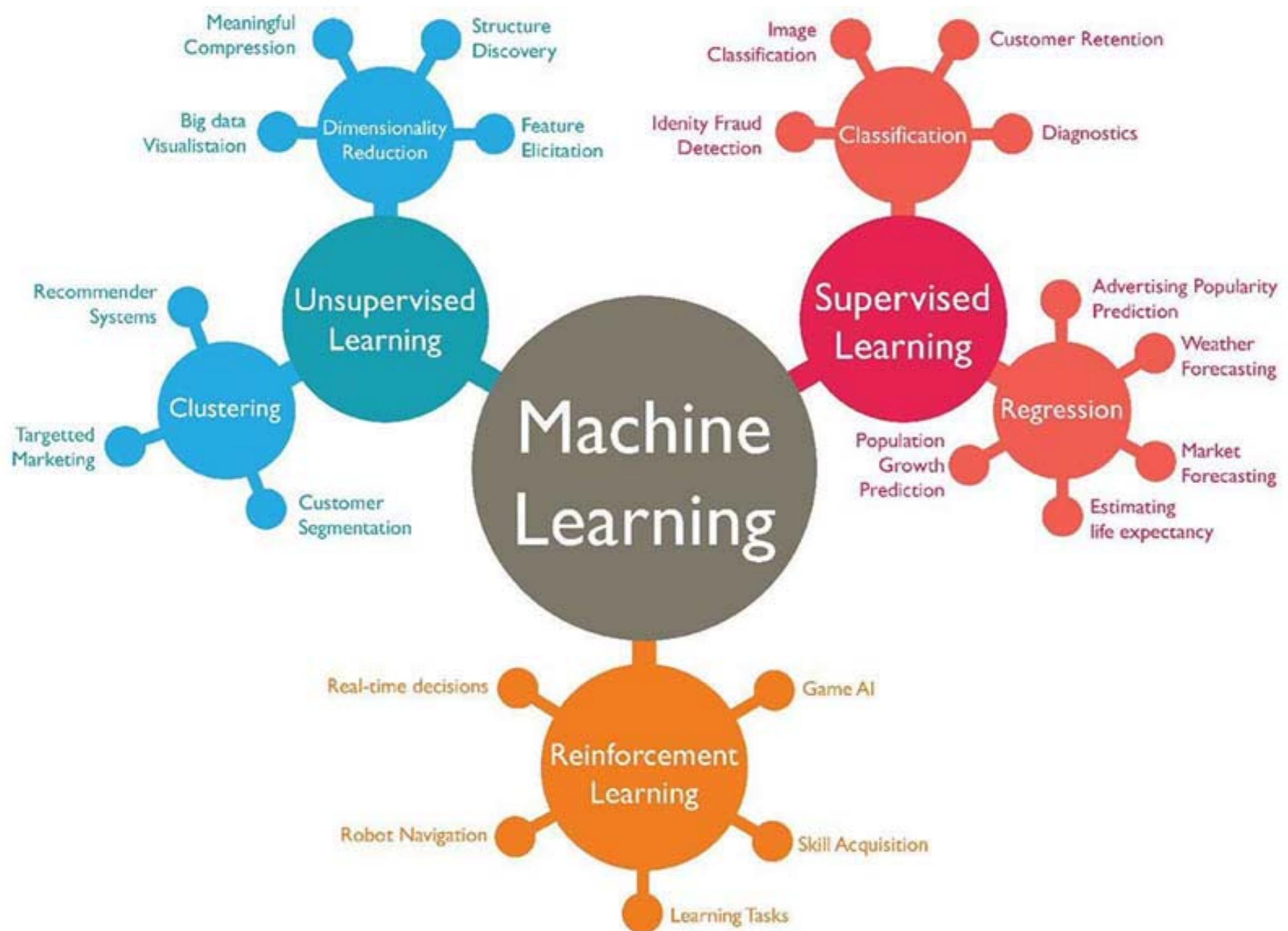


Image source : www.techleer.com

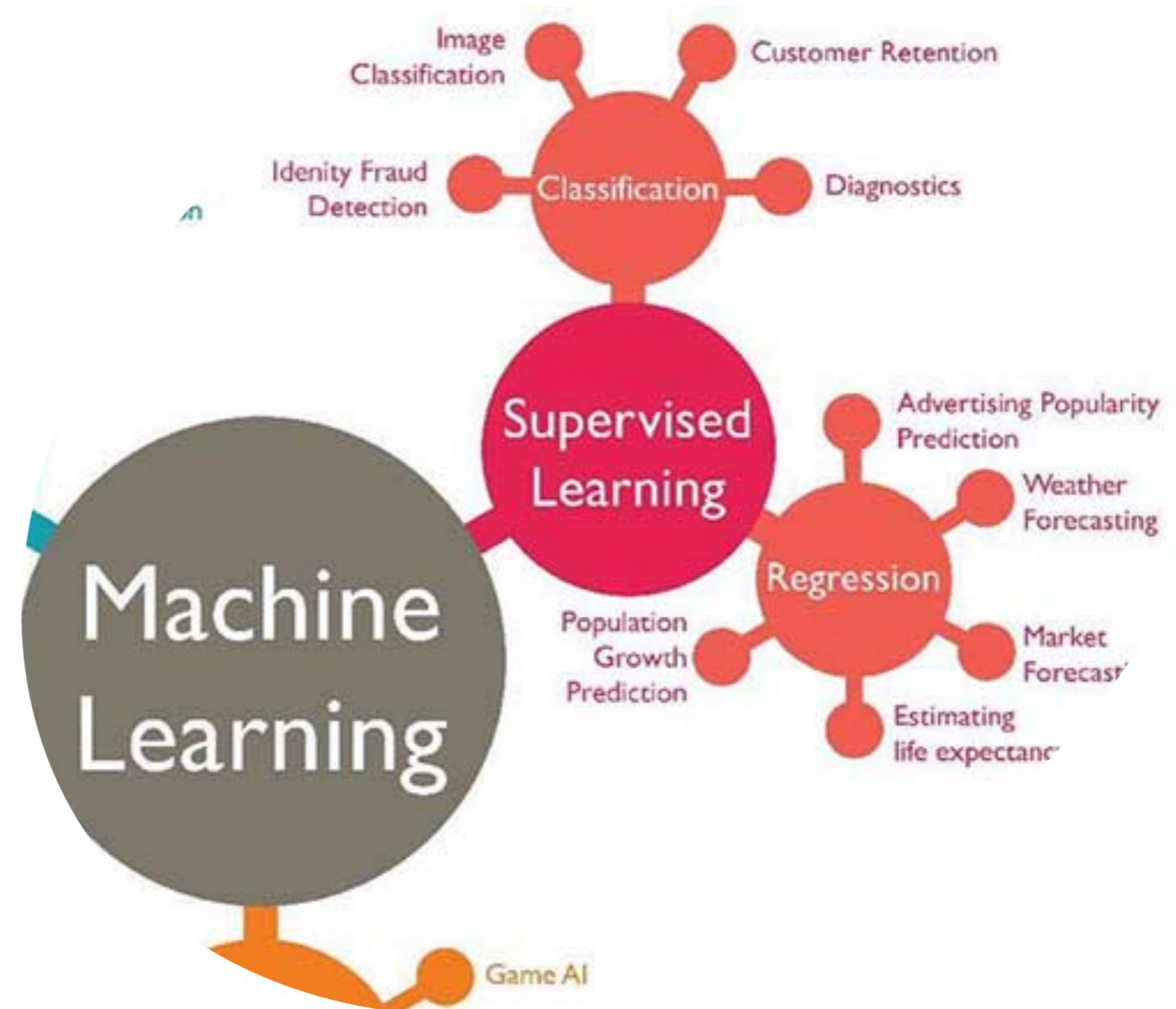


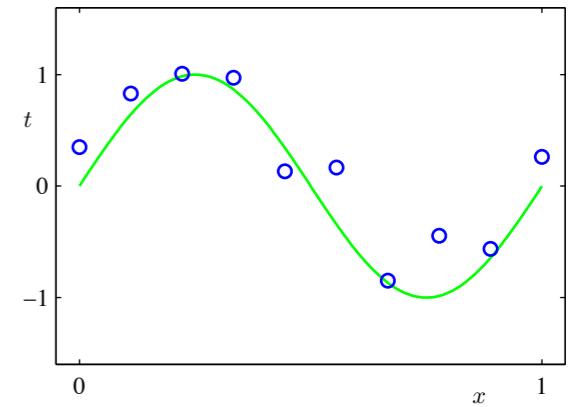
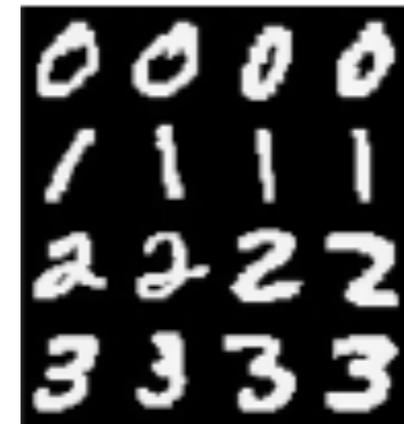
Image source : www.techleer.com

Supervised learning

Associative learning



Dataset



features: $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$

$\mathbf{x} = 2$

$\mathbf{x} = 0.25$

targets: $\{t_1, \dots, t_N\}$

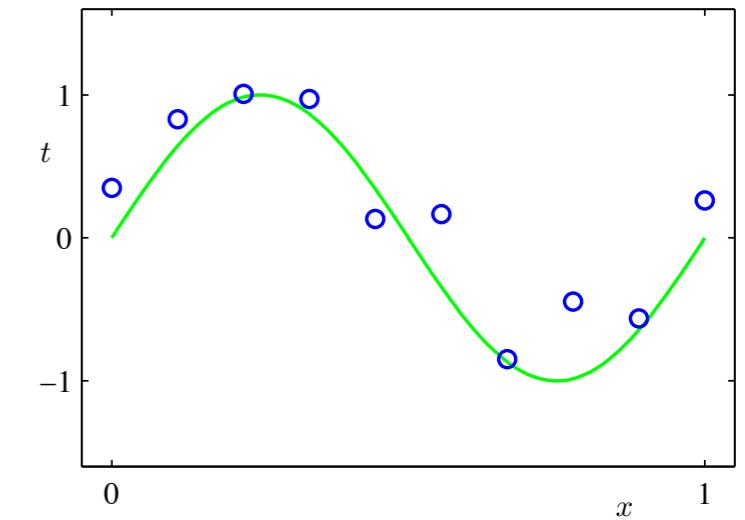
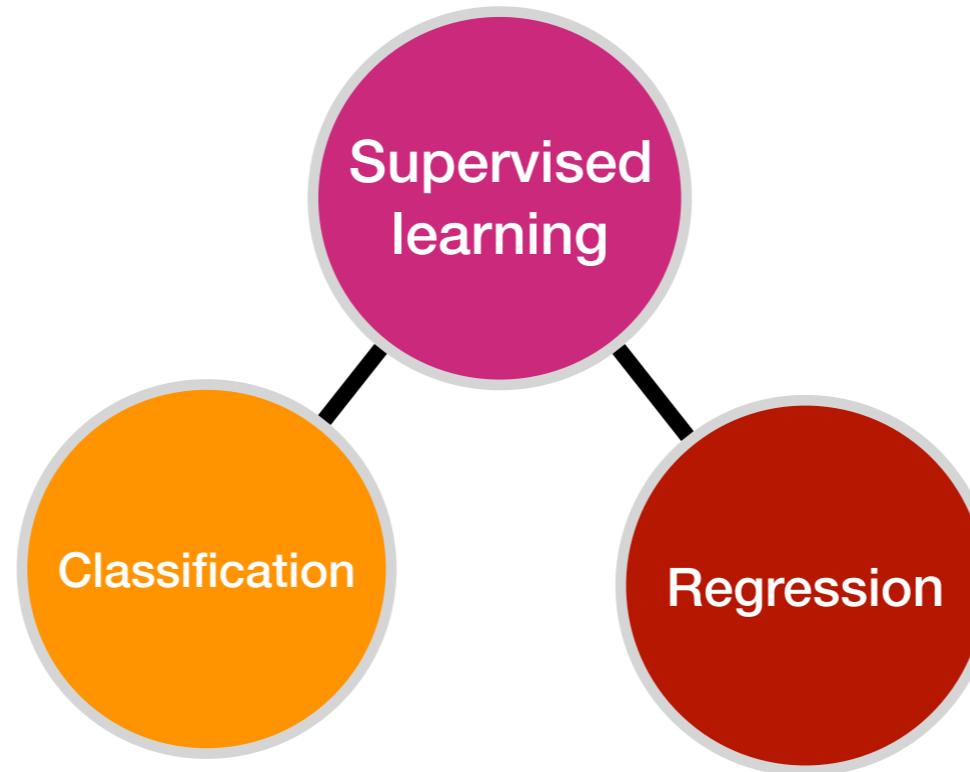
$t = 2$

$t = 0.707$

discrete

continuous

Supervised learning

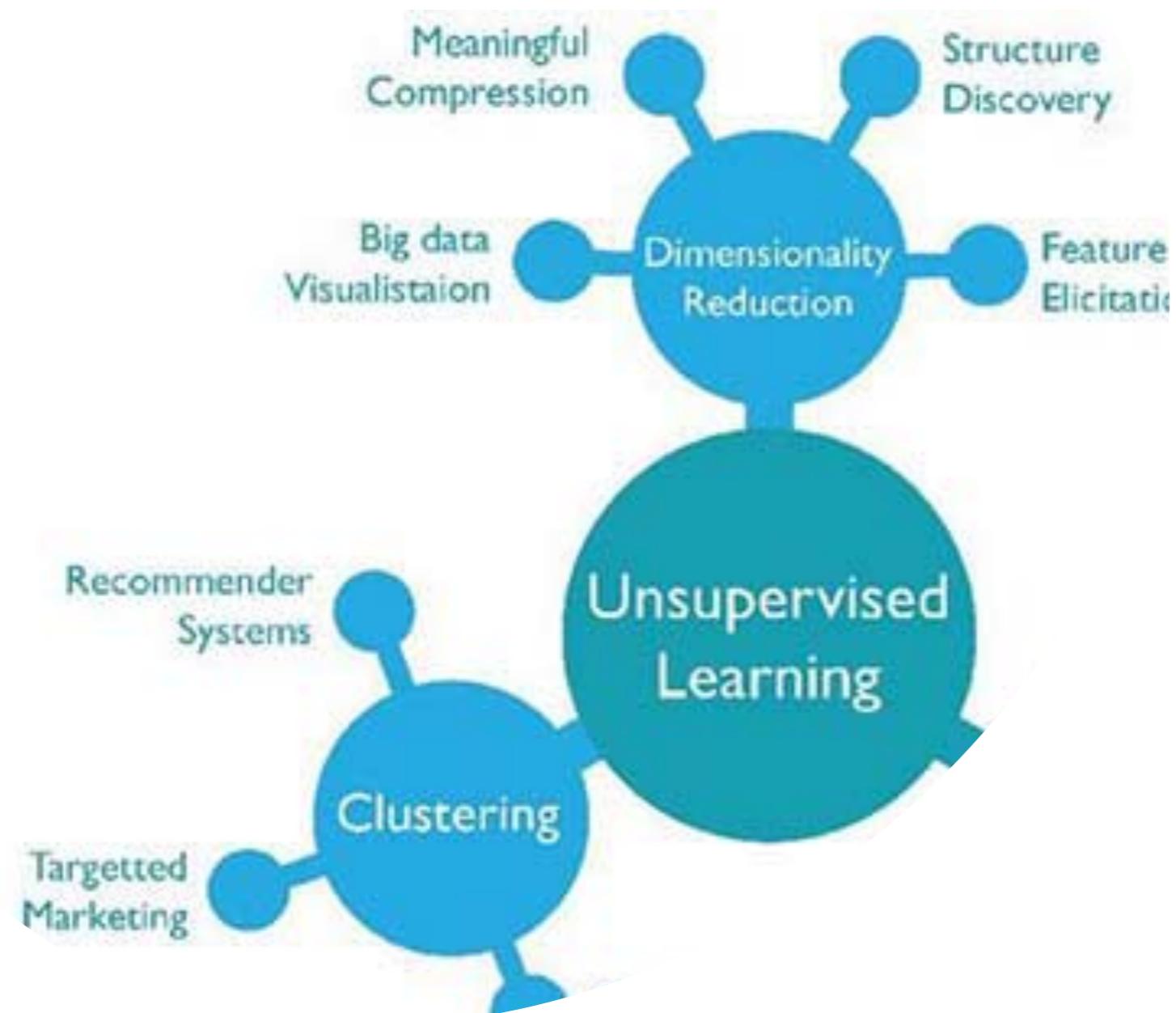


Task: Find function f such that $f(x) \approx t$ for all known
and unknown (x, t)

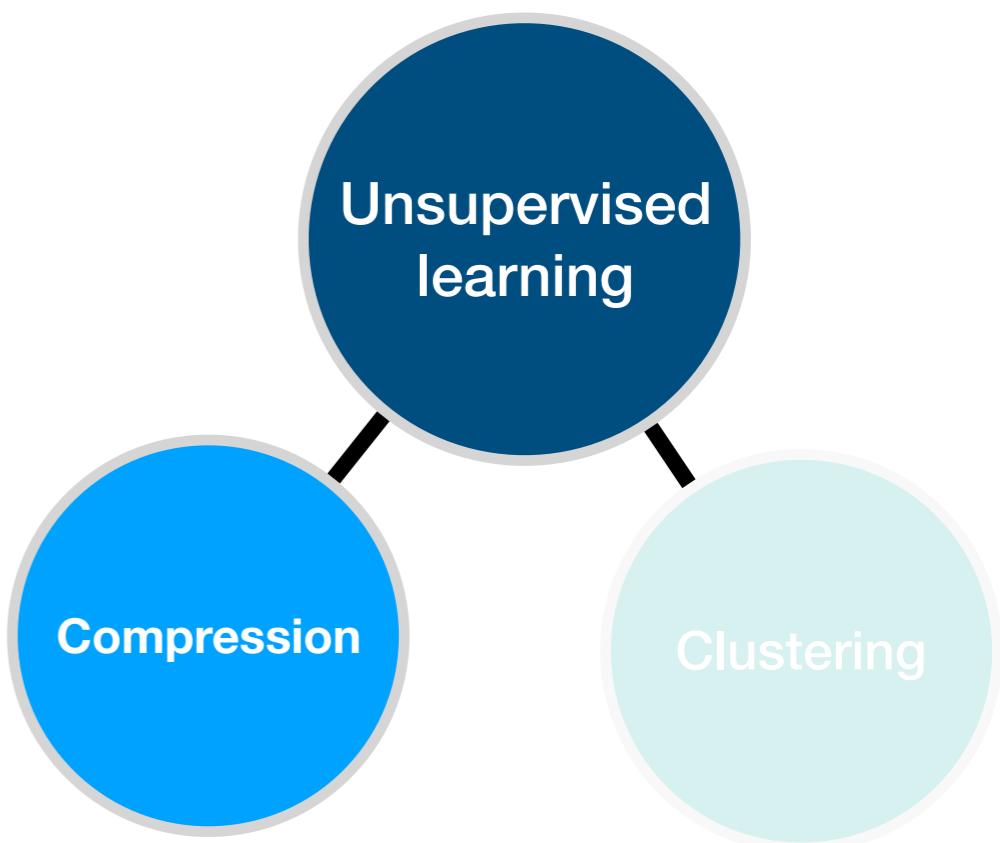
validation / test data

training pair

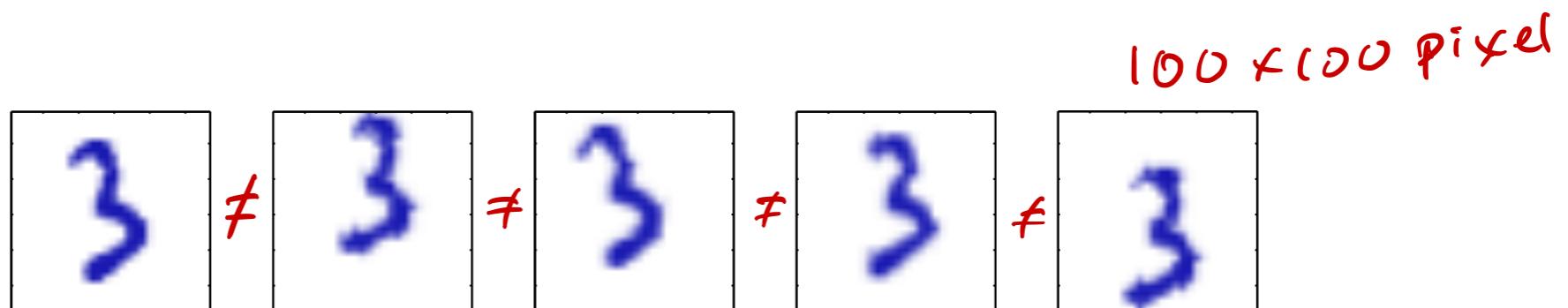
training
data



Unsupervised learning



Dataset:

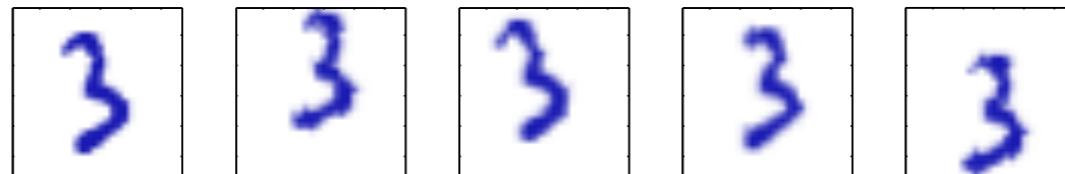


Task: compress images

why? reduces required storage

Unsupervised learning

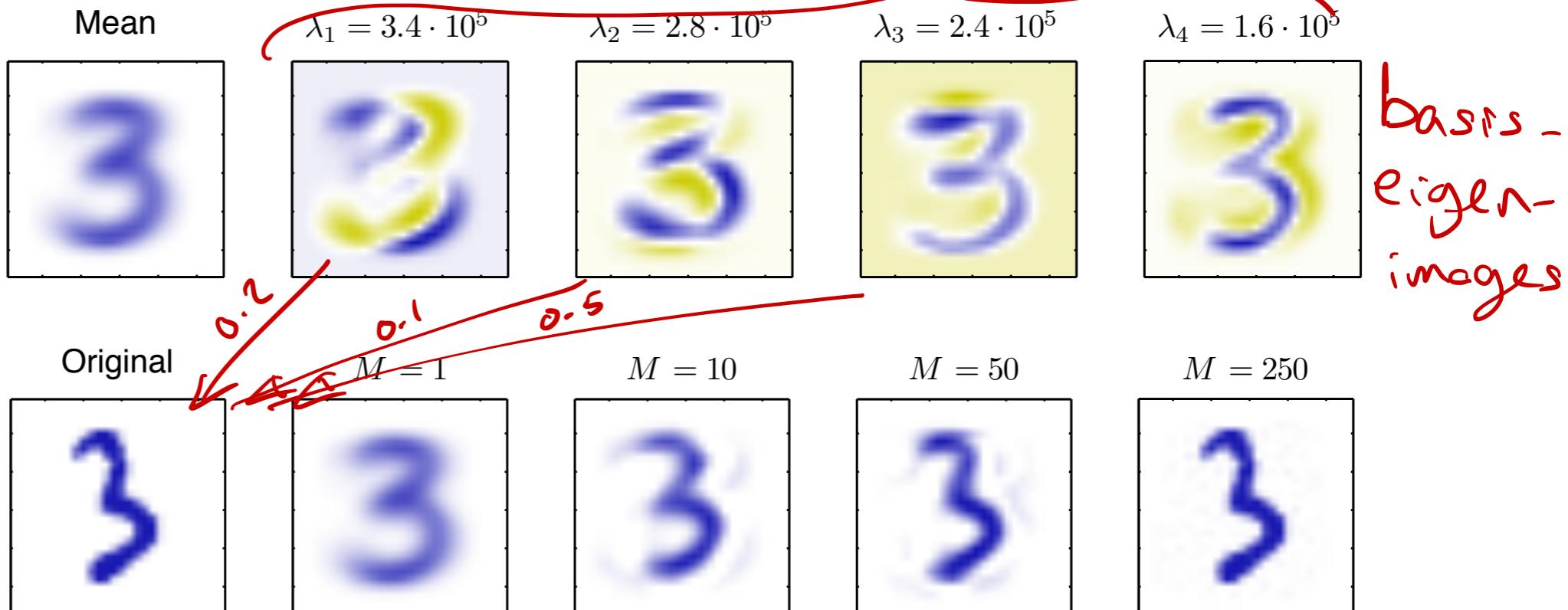
Dataset:



Task: Compress image

Method: expand along principle components (PCA)

need to save
M principle
components
(size 100x100)

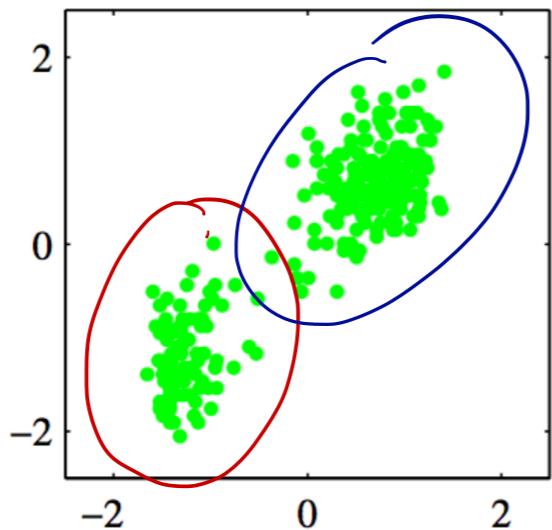


Result:

per image: only
need to save M
IR - numbers

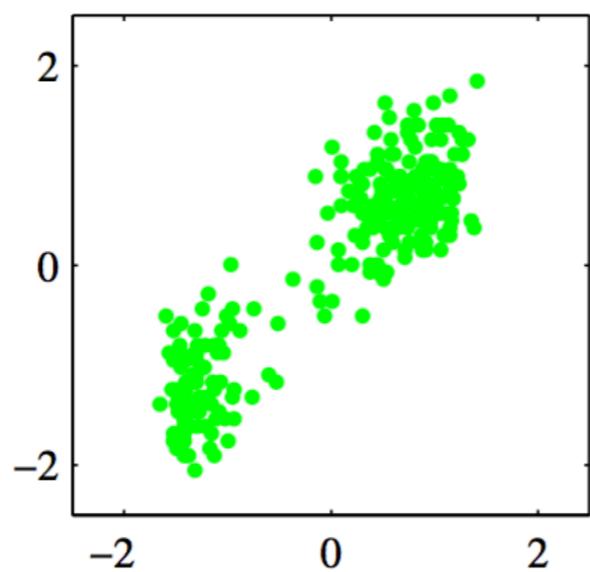
Unsupervised learning

Dataset:

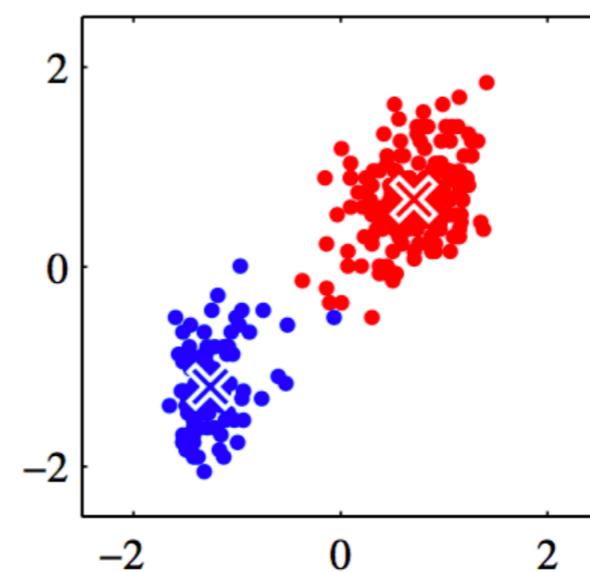


Task: Assign every datapoint to a cluster (hidden class variable)

Result:



Dataset



Final clustering

Other types of learning

Semi-supervised learning

- data points: $\{x_1, \dots, x_n\}$
 $\text{"cat"} \quad \text{"dog"}$
- targets: $\{t_1, t_2, \dots, t_k\} (k < n)$
- Not all datapoints have a known target/label!
only the first $k < n$
- Use all data, also those with unknown target, to create predictor.
image classification

Other types of learning

Reinforcement Learning

- Dynamic environment: provides information on its state. *S*
- Agent: takes actions, receives rewards from environment. *a*
- Task: maximize total reward *r*
- Learning by trial and error
- Application: *games*



Overview

1. Overview

1. What is machine learning?
2. Different types of machine learning

2. Recap: probability theory

Probability theory

Probability theory (Bishop)

Provides a consistent framework for the quantification and manipulation of uncertainty.

everything I don't know

Uncertainty in pattern recognition

- Noise on measurements. *unknown noise sources*
- Finite size datasets. *the fewer samples available
the less you know about the
true distribution*

Probability theory

Frequentist interpretation

- Probability of event: fraction of times event occurs in experiment

Bayesian approach

- Probability: quantification of plausibility or the strength of the belief of an event.

think science

think intuitive

Random variables

Random variable $X \sim p(x)$ *n measurable*

- Stochastic variable sampled from a set of possible outcomes

$$p(x) \geq 0 \quad \forall x \in \mathcal{X}$$

- Discrete or continuous

- Probability distribution $p(X)$, $p_{X=x}$, $P(X = x)$, $p(x)$
capital lowercase

Examples of discrete random variables:

- Throwing a die: $X = \# \text{ of eyes}$ $\mathcal{X} = \{1, 2, \dots, 6\}$

- Flipping a coin: $X = \text{side of the coin}$
 $\mathcal{X} = \{\text{heads}, \text{tails}\}$

$$\sum_{x \in \mathcal{X}} p(x) = 1$$

Two discrete random variables (I)

2 random variables

$X \neq Y$

N trials: sample both X and Y .

Joint probability

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

Marginal probability of X : $p(X = x_i) = c_i/N$

$$c_i = \sum_{j=1}^3 n_{ij} \quad n_{ij} = p(X = x_i, Y = y_j) \cdot N$$

$$p(X = x_i) = \sum_{j=1}^3 p(X = x_i, Y = y_j)$$

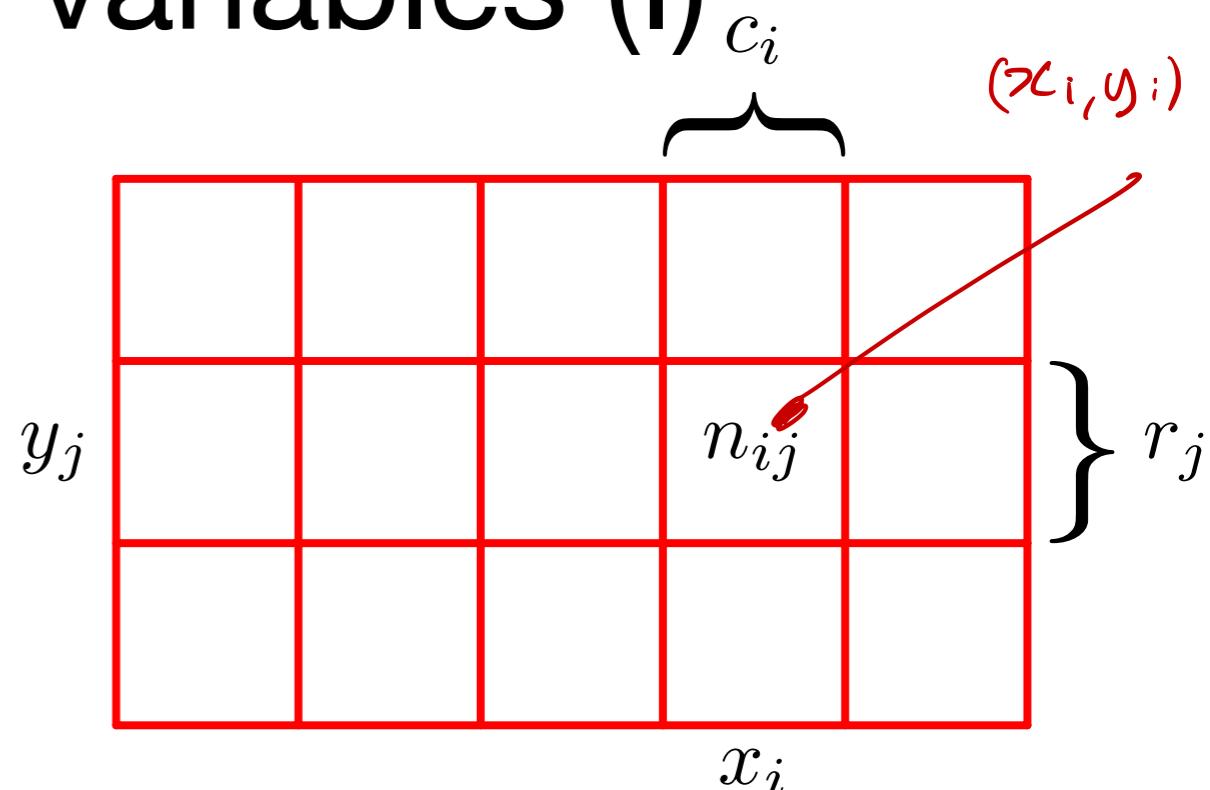


Figure: 2 random variables (Bishop 1.10)

Sum rule /
marginalisation
"projection"

Two discrete random variables (II)

- 2 random variables

- Conditional probability of Y given X:

$$P(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i} \leftarrow \text{renorm.}$$

- Remember: $p(X = x_i) = \frac{c_i}{N}$

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = \frac{P(Y = y_j | X = x_i) \cdot c_i}{N}$$

$$P(X = x_i, Y = y_j) = P(Y = y_j | X = x_i) P(X = x_i)$$

Product rule

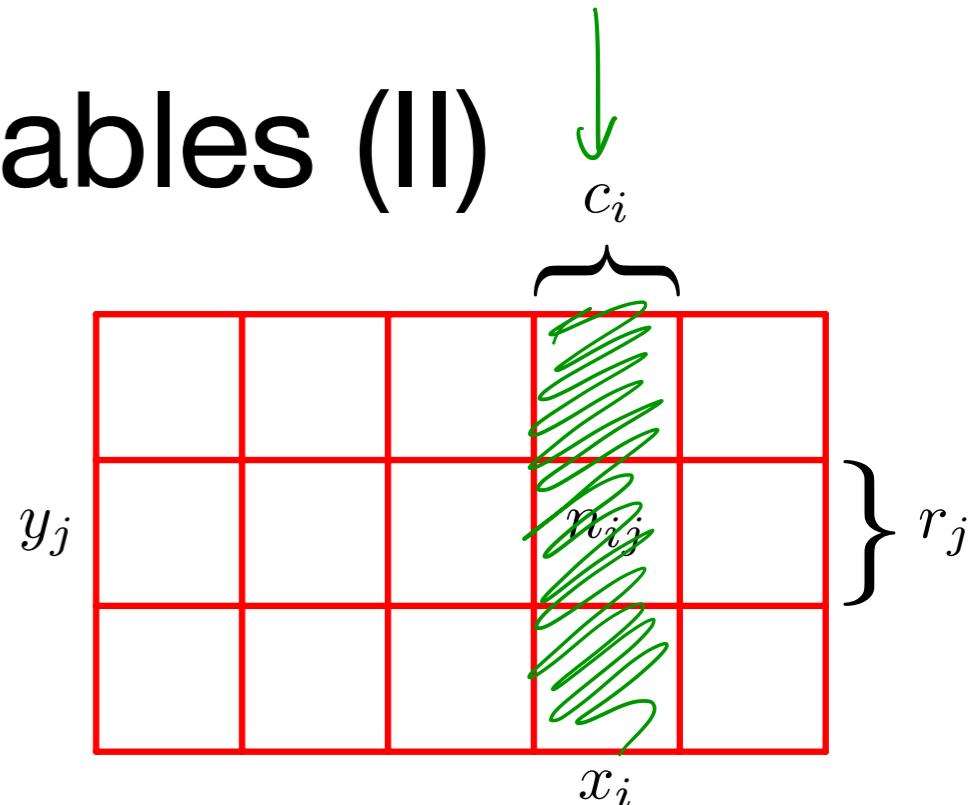


Figure: 2 random variables (Bishop 1.10)

Example: Marginal & Conditional distributions

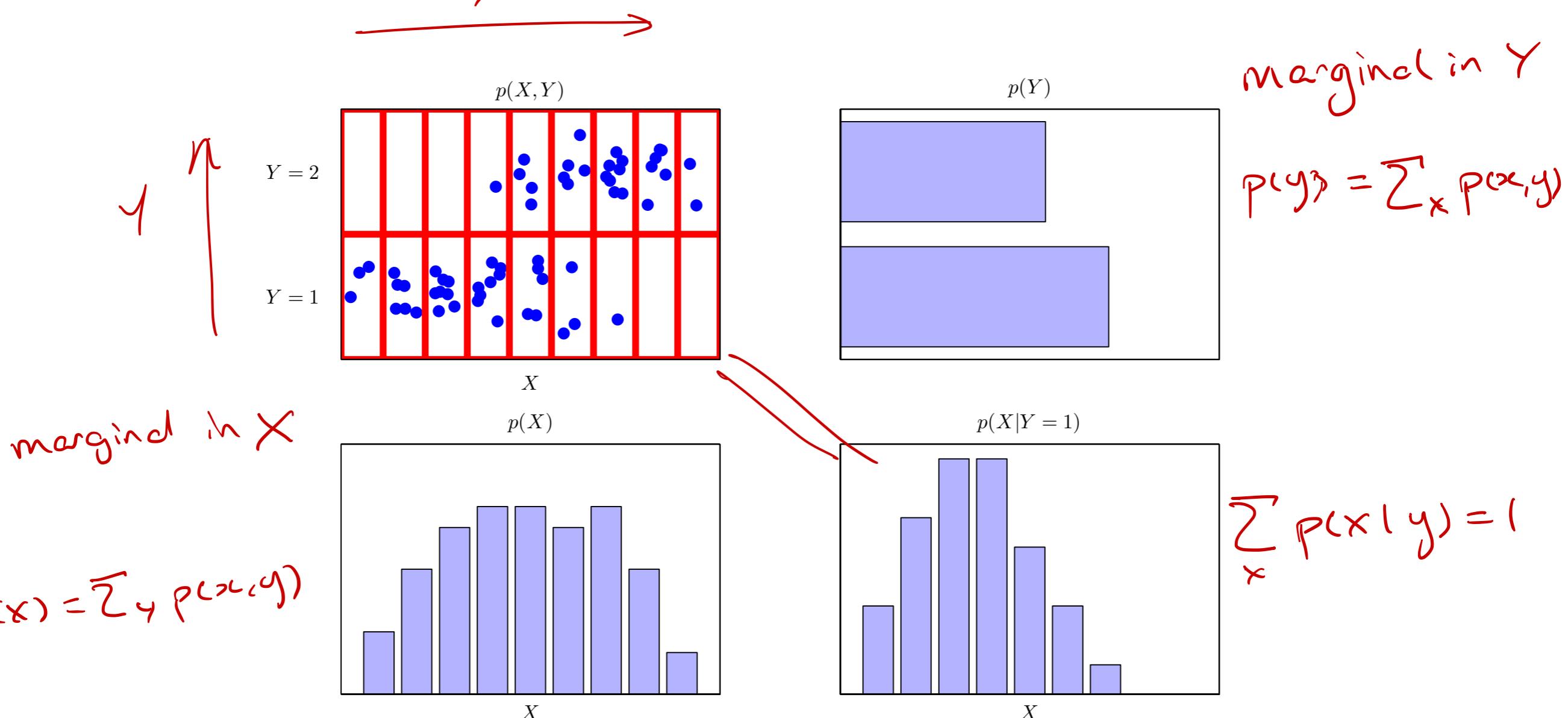


Figure: Marginal and conditional distributions (Bishop 1.11)

Continuous Random Variables

- Probability of $x \in \mathbb{R}$ falling in the interval $[x, x + dx]$ is given by $|p(x)dx|$ prob. mess
- $p(x)$: probability density over x
- Probability over finite interval $p(x \in (a, b)) = \int_a^b dx p(x)$, $b > a$
- Positivity: $p(x) \geq 0$
- Normalization: $\int_{-\infty}^{\infty} dx p(x) = 1$
- Change of variables $x = g(y)$, probabilities in $(x, x + dx)$ must be transformed to $(y, y + dy)$
$$|p_x(x)dx| = |p_y(y)dy| \rightarrow p_y(y) = p_x(g(y)) \left| \frac{dx}{dy} \right|$$

Continuous Random Variables

$$F(x) = P(X \leq x)$$
$$= \int_{-\infty}^x dx' p(x')$$

$$p(x) = \frac{dF}{dx}(x)$$

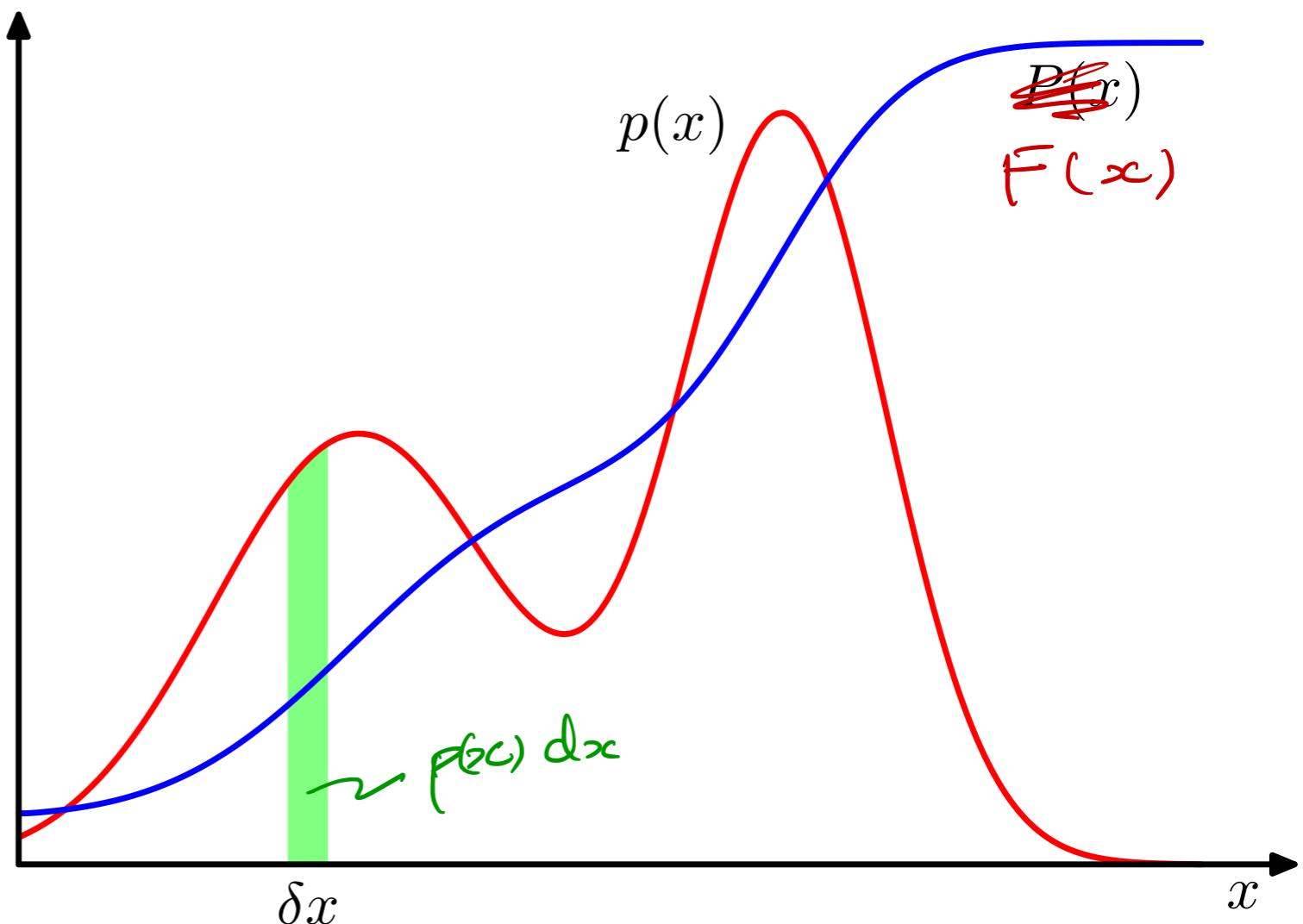


Figure: probability density and cumulative distribution function (Bishop 1.12)

The Rules of Probability Theory

For random variables $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$:

$$A = (a, b) \quad p(X \in A) = \int_a^b dx \cdot p(x)$$

| | Discrete | Continuous |
|---------------|-------------------------------------------|---------------------------------------------|
| Additivity | $p(X \in A) = \sum_{x \in A} p(x)$ | $p(X \in A) = \int_{x \in A} dx \cdot p(x)$ |
| Positivity | $p(x) \geq 0$ | $p(x) \geq 0$ |
| Normalization | $\sum_x p(x) = 1$ | $\int_{\mathcal{X}} p(x) dx = 1$ |
| Sum Rule | $p(x) = \sum_{y \in \mathcal{Y}} p(x, y)$ | $p(x) = \int_y dy \cdot p(x, y)$ |
| Product Rule | $p(x, y) = p(x y)p(y)$ | $p(x, y) = p(x y)p(y)$ |

Bayes' Theorem

- ▶ Product rule $p(x, y) = p(x|y)p(y)$
- ▶ Symmetry property $= p(y|x)p(x)$
 ↑
 likelihood
 prior
- ▶ Bayes rule $\rightarrow p(y|x) = \frac{p(x|y)p(y)}{p(x)}$
 ↑
 marginal
 likelihood
 evidence / normalization
 constant
- ▶ Denominator: $p(x) = \sum_y p(x|y)p(y)$
 $\Rightarrow 1 = \sum_y \left(\frac{p(x|y)p(y)}{p(x)} \right)$

Bayes Theorem

Bayes rule

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

- $p(y)$: the prior probability of $Y = y$
- $p(y | x)$: the posterior probability of $Y = y$
- $p(x | y)$: the likelihood of $Y = y$ given $X = x$
- $p(x)$: the evidence for $X = x$

Example: Fruit in Boxes

- ▶ Random variables:

- ▶ Prior Box distribution:

$$p(B = r) =$$

$$p(B = b) =$$

- ▶ Conditional probabilities of Fruit given Box

- ▶ Marginal Fruit distributions:

$$p(F = a) =$$

$$p(F = o) =$$

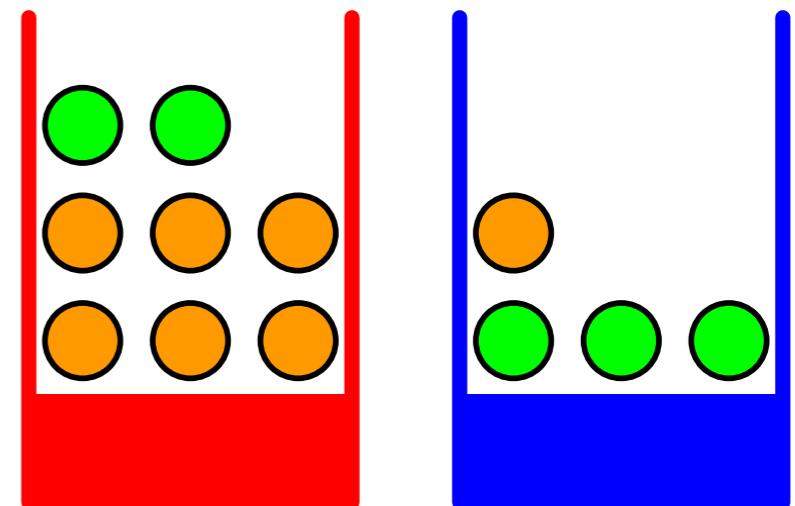
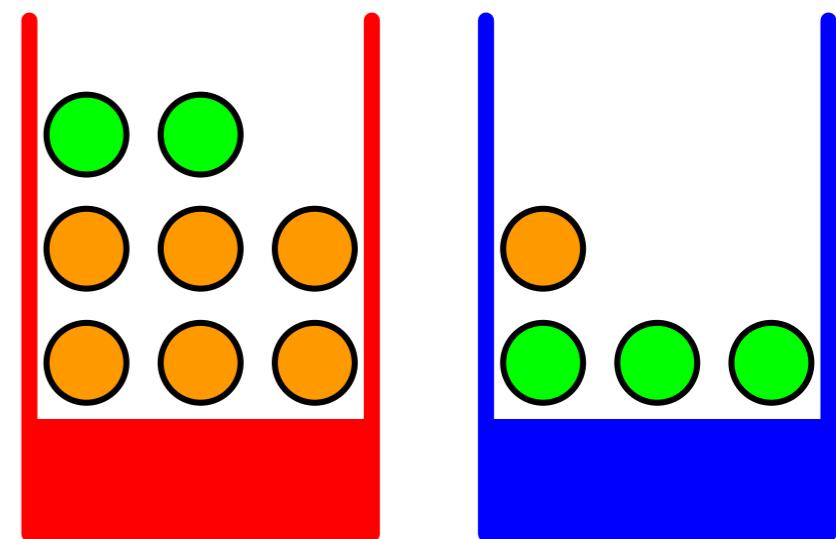


Figure: coloured boxes containing apples and oranges (Bishop 1.9)

Example: Fruit in Boxes

- ▶ Prior: $p(B = r) = 4/10$ & $p(B = b) = 6/10$
- ▶ Marginal: $p(F=a) = 11/20$ & $p(F=o) = 9/20$
- ▶ Posterior probability of Box color given observed fruit

$$p(B = r|F = o) =$$



- ▶ prior probability of red box:

$$p(B = r) = 4/10$$

- ▶ After observing an orange the probability of observing a red box is now larger than observing a blue box!

Figure: coloured boxes containing apples and oranges (Bishop 1.9)

Independent Random Variables

Two random variables X and Y are *independent* iff measuring X gives no information on Y , and vice versa.

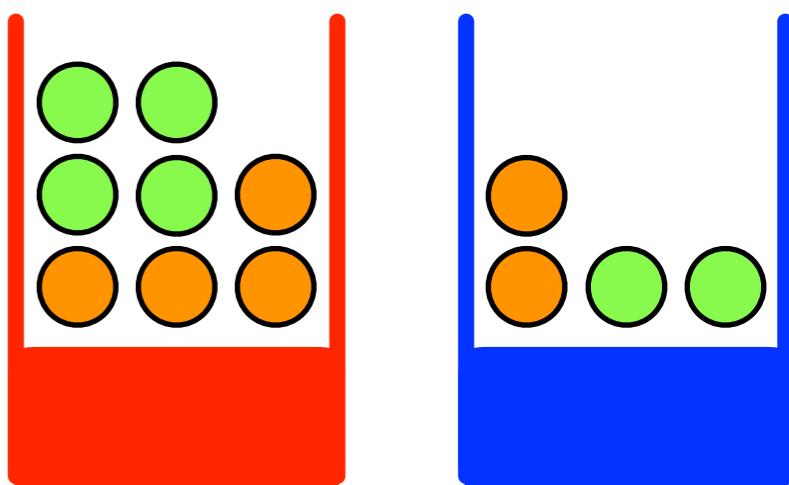
- Formally: X and Y are called independent if

for all $x \in \mathcal{X}, y \in \mathcal{Y}$

- Equivalent to

$$p(x|y) =$$

- Example:



Expectations

- random variable $X \in \mathcal{X}$ and function $f : \mathcal{X} \rightarrow \mathbb{R}$

$$\mathbb{E}[f] = \mathbb{E}_{x \sim p(X)}[f(x)] =$$

- For N points drawn from $p(\mathbf{X})$:

$$\mathbb{E}[f] =$$

- Conditional expectation:

$$\mathbb{E}[f | y] = \mathbb{E}_{x \sim p(X|Y=y)}[f(x)]$$

Variance

- The expected quadratic distance between f and its mean $\mathbb{E}[f]$

$$\text{var}[f] =$$

Covariance between 2 random variables

- ▶ Measures the extent to which X and Y vary together

$$\text{cov}[x, y] =$$

- ▶ Vectors of random variables \mathbf{x} and \mathbf{y} , covariance matrix:

$$\text{cov}[\mathbf{x}, \mathbf{y}] =$$

- ▶ Independent variables:

$$\mathbb{E}[xy] =$$

- ▶ Note: $\text{cov}[x, y] = 0$ does not imply

- ▶ $\text{cov}[\mathbf{x}] \equiv \text{cov}[\mathbf{x}, \mathbf{x}]$

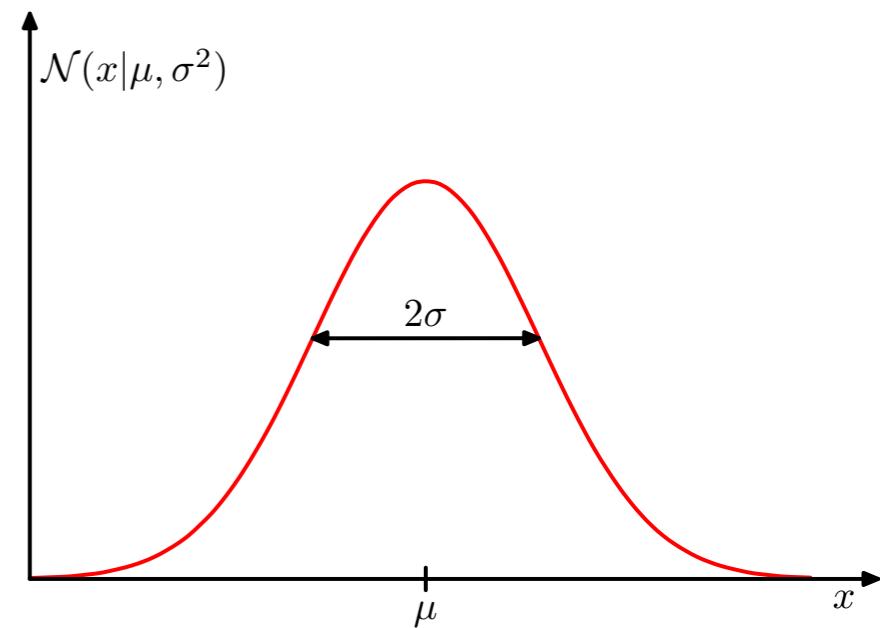
Gaussian Distribution

- Real valued stochastic variable X

$$\mathcal{N}(x|\mu, \sigma^2) =$$

- Mean:

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x \, dx =$$



$$\int_{-\infty}^{\infty} e^{-x^2} dx = \sqrt{\pi}$$

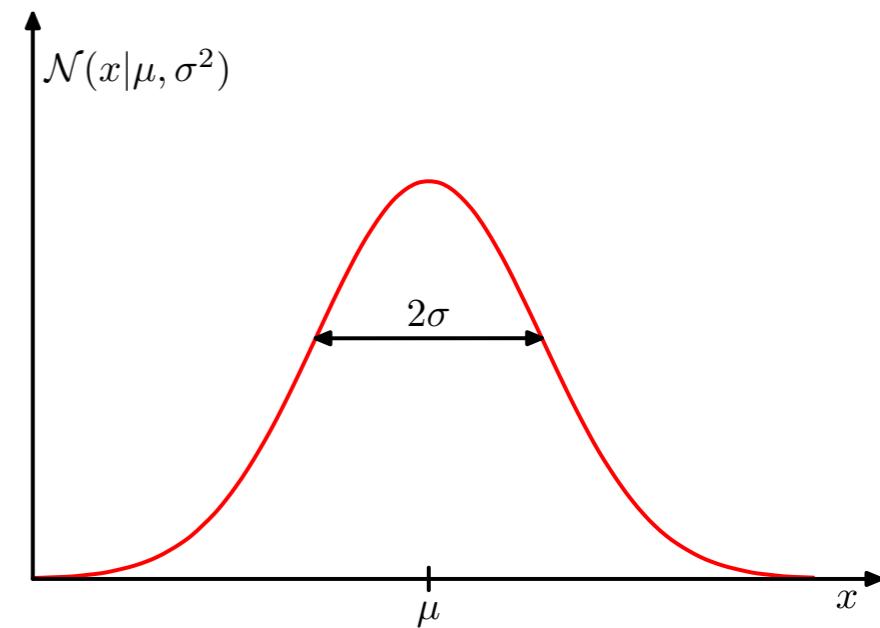
Gaussian Distribution

- Real valued stochastic variable X

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2} (x - \mu)^2\right\}$$

- Variance: $\text{var}[x] = \mathbb{E}[(x - \mathbb{E}[x])^2]$

$$\text{var}[x] =$$



$$\int_{-\infty}^{\infty} e^{-x^2} dx = \sqrt{\pi}$$

$$\int_{-\infty}^{\infty} x^2 e^{-ax^2} dx = -\frac{\partial}{\partial a} \int_{-\infty}^{\infty} e^{-ax^2} dx =$$

Multivariate Gaussian Distribution

- D -dimensional vector $\mathbf{x} = (x_1, x_2, \dots, x_D)^T$

- $\mathcal{N}(\mathbf{x}|\mu, \Sigma) =$

$$|\Sigma| = \det \Sigma$$

- $\Sigma =$ (D x D matrix)

- $\mathbb{E}[\mathbf{x}] =$

$$\int \exp\left\{-\frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x}\right\} d^D x = \frac{(2\pi)^{D/2}}{|\mathbf{A}|^{1/2}}$$