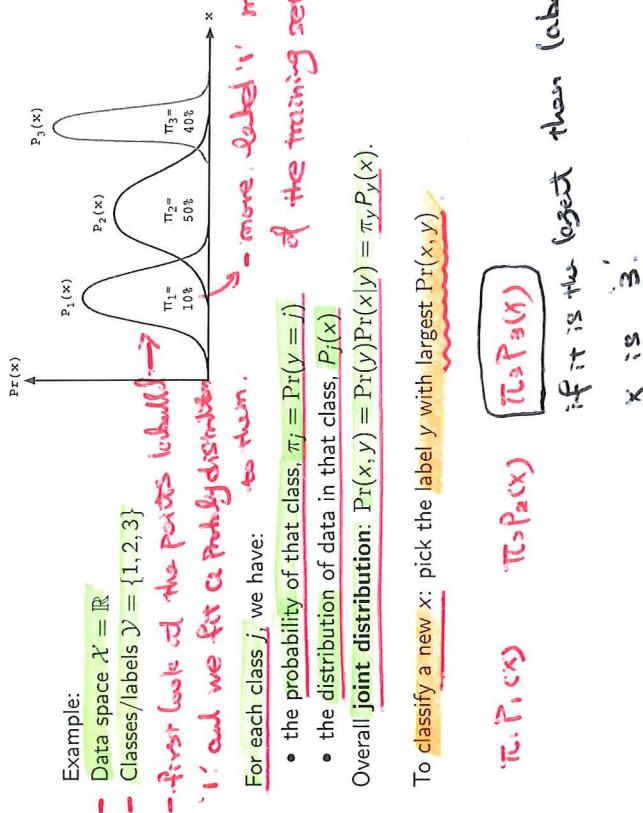


Generative models

WEEK 2-1 1st (Two) days:

The generative approach to classification

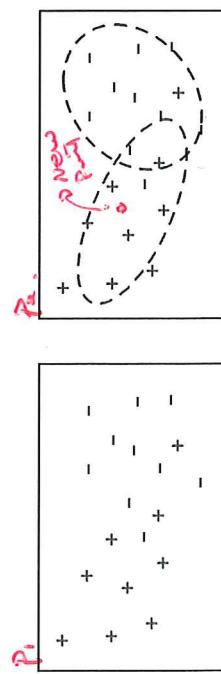
- Generative approach based on probability distributions.
- Main idea & Generative approach:
 - fit each class separately with a probability distribution.



$$\pi, \pi_1(x), \pi_2(x)$$

$$\pi, p_1(x)$$

The generative approach to classification



- Training set with 15-20 points
- 2 labels, plus 1 minus
- look at classes labeled and we fit a model to them

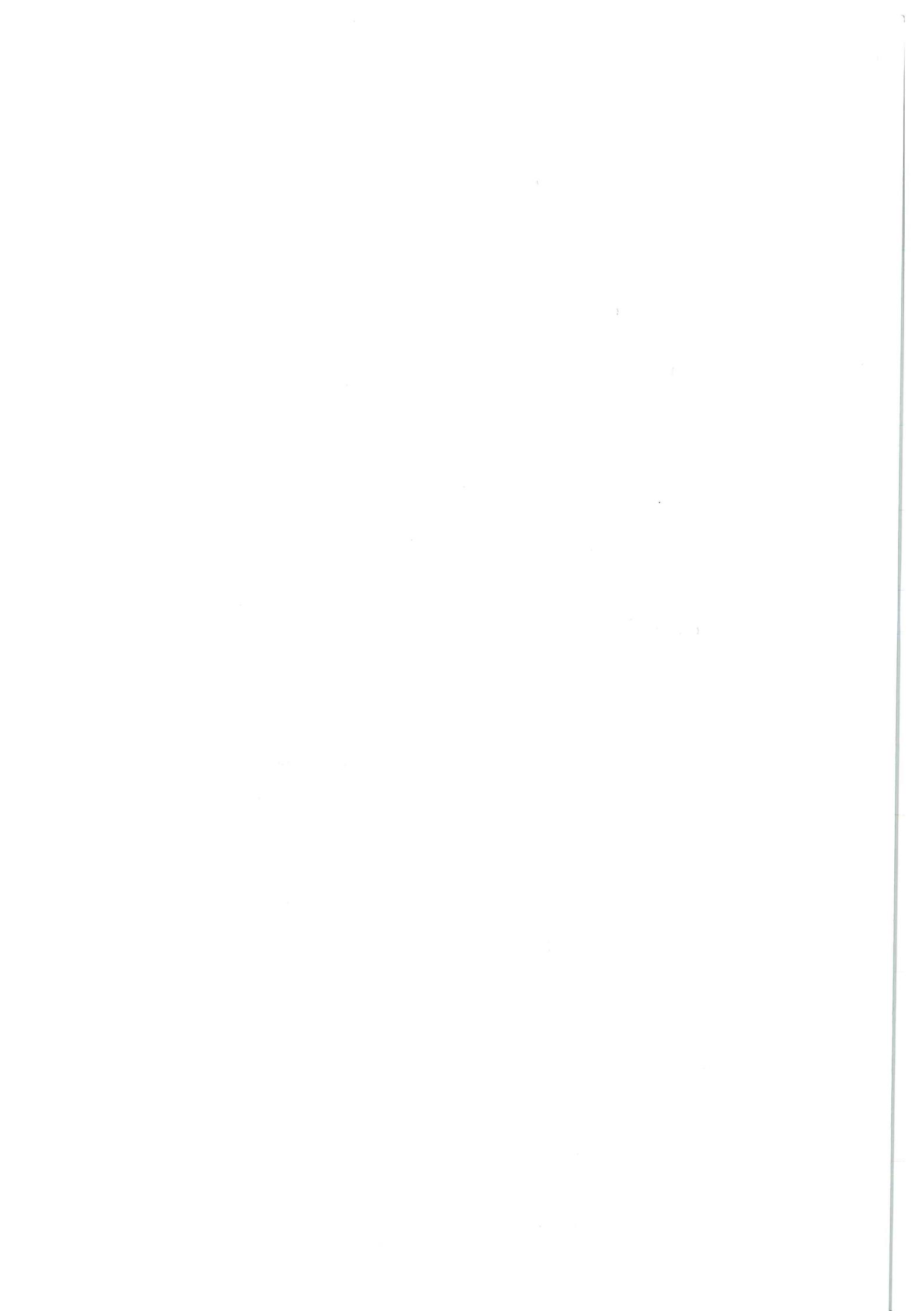
- then look at minus we ...
- ellipse shaped distribution that fit into plus & minus.

The learning process:

- Fit a probability distribution to each class, individually

To classify a new point:

- Which of these distributions was it most likely to have come from?



WEEK 2 - 2

14/ Tues(2nd)

Probability review I:

Probability spaces, events, and conditioning

Probability spaces

You roll two dice.

Q: What is the probability they add to 10?

The probability space has two components:
① Sample space (space of outcomes). **the set of all possible outcomes.**

$$\Omega = \{(1,1), (1,2), \dots, (1,6), (2,1), \dots, (6,6)\}$$

$$= \{1, 2, \dots, 6\} \times \{1, 2, \dots, 6\} = \{(1,2,3,4,5,6)\}^2 = 36.$$

- ② Probabilities of outcomes, summing to 1. Probability of each of these outcomes.
- the number of possible outcomes = $6 \times 6 = 36$.
- Each outcome has prob $\frac{1}{36}$.
↳ Probability space.

Topics we'll cover

Events

Probability space:

- Outcomes: $\Omega = \{\text{all possible pairs of dice rolls}\}$
- Every pair $z = (z_1, z_2) \in \Omega$ has probability $1/36$.

Event of interest: the two dice add up to 10.

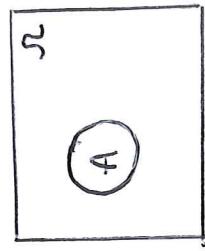
$A \subseteq \Omega$ (An Event is just a subset of the possible outcomes).

$A = \{(z_1, z_2) : z_1 + z_2 = 10\}$

$$= \{(4,6), (5,5), (6,4)\}$$

$$P(A) = 3 \times \frac{1}{36} = \frac{1}{12}$$

↓
3 outcomes
Probability $1/12$.



Multiple events

- You have ten coins. Nine are fair, but one is a bad coin that always comes up tails.
- You close your eyes and pick a coin at random.
 - You toss it four times, and it comes up tails every time.
- What is the probability you picked the bad coin?

$$\text{Sample space} = \{H, T\}^4 = \{HHTT, HTTH, THHT, TTHH, HHTH, HTTH, THHT, TTHH, HHHT, HHTT\}$$

$$= 10^4 = 10 \times 10 \times 10 \times 10 = 10000$$

- Ten coins: nine are fair, one is a bad coin that always comes up tails.
- You pick a coin at random, toss it four times, and it's tails every time.

$A = \text{Picked the bad coin}$

$$= \{\text{bad coin}\}$$

$$B = \text{all coins are tails}$$

$$= \{(H, H, H, H), (T, T, T, T)\}$$

$$\Pr(A \cap B) = \Pr(\text{bad coin}) \Pr(\text{all tails | bad coin})$$

$$= \frac{1}{10} \times \frac{1}{10} = \frac{1}{100}$$

Conditioning

You have ten coins. Nine are fair, but one is a bad coin that always comes up tails.

- You close your eyes and pick a coin at random.
 - You toss it four times, and it comes up tails every time.
- What is the probability you picked the bad coin?

to come, 4 times.

Conditioning formula: $\Pr(A \cap B) = \Pr(A) \Pr(B|A)$

In our example:

- A: the bad coin is chosen
- B: all four tosses are tails

Want $\Pr(A|B)$

$$\Pr(A|B) = \frac{\Pr(A \cap B)}{\Pr(B)} = \frac{\frac{1}{10}}{\frac{1}{100}} = 0.1$$

- Ten coins: nine are fair, one is a bad coin that always comes up tails.
- You pick a coin at random, toss it four times, and it's tails every time.

Event A: the bad coin is chosen. Event B: all tails

$$\Pr(A|B) = \frac{\Pr(A \cap B)}{\Pr(B)} = \frac{\Pr(\text{bad coin | all tails})}{\Pr(\text{all tails})}$$

$$= \frac{1}{10} \times 1 + \frac{9}{10} \times \left(\frac{1}{2}\right)^4 = \frac{5}{32}$$

$$\Pr(A|B) = \frac{\Pr(A \cap B)}{\Pr(B)}$$

$$= \frac{\frac{1}{10}}{\frac{5}{32}} = 0.64$$

Bayes' rule

Two events A, B

- We are interested in A
- We can observe B

If we find out B occurred, how does it alter the probability of A ?

$$\boxed{\text{Bayes' rule: } \Pr(A|B) = \Pr(A) \times \frac{\Pr(B|A)}{\Pr(B)}}$$

Corrected factor

$$\Pr(A|B) = \frac{\Pr(\text{AND})}{\Pr(B)} = \frac{\Pr(A) \Pr(B|A)}{\Pr(B)}$$

Tell. If a coin is flipped 5 times, what is the

size of the sample space for the experiment?

- the sample space of a fair coin flip is $\{H, T\}$.
- the sample space of a coin flips is all 2^5 possible sequences of outcomes $2^5 = 32$.

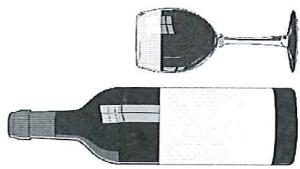


A classification problem

You have a bottle of wine whose label is missing.

Week-3. 15 Jan 2018.

Generative modeling in one dimension



Which winery is it from, 1, 2, or 3?

Solve this problem using visual and chemical features of the wine.

Topics we'll cover

The data set

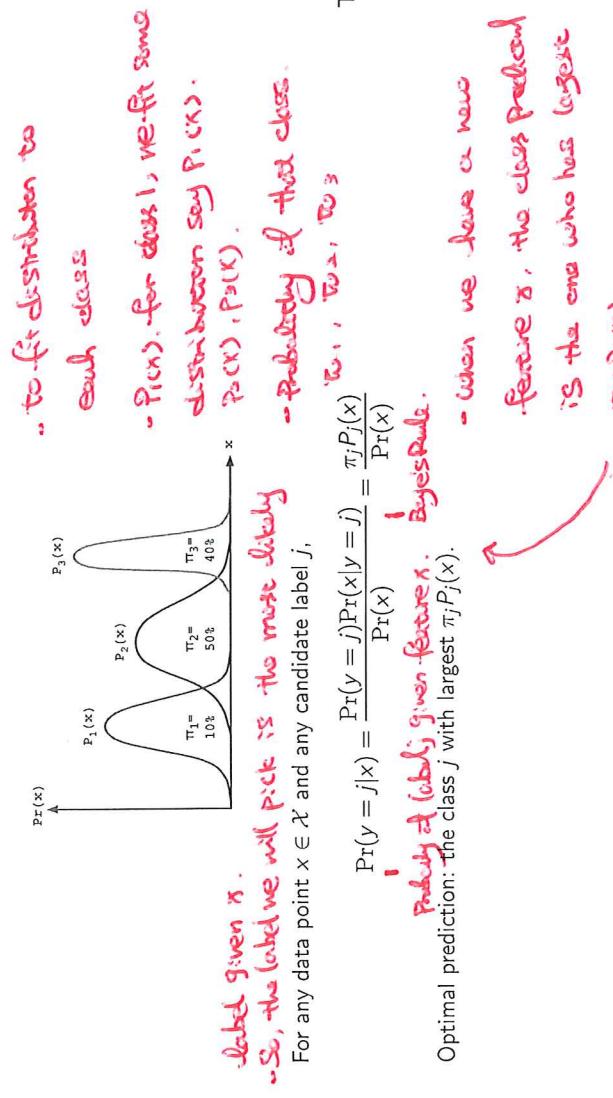
Training set obtained from 130 bottles *Train Set*.

- Winery 1: 43 bottles
- Winery 2: 51 bottles
- Winery 3: 36 bottles
- For each bottle, 13 features: *13 dimensional vectors*.
 - 'Alcohol', 'Malic acid', 'Ash', 'Alcalinity of ash', 'Magnesium', 'Total phenols', 'Flavanoids', 'Nonflavanoid phenols', 'Proanthocyanins', 'Color intensity', 'Hue', 'OD280/OD315 of diluted wines', 'Proline'

Also, a separate test set of 48 labeled points. *Test Set*.

- *use this data to build classifier that takes the features from a new bottle and predict its label.*

Recall: the generative approach



Fitting a generative model

Training set of 130 bottles:

- Winery 1: 43 bottles, winery 2: 51 bottles, winery 3: 36 bottles
- For each bottle, 13 features: 'Alcohol', 'Malic acid', 'Ash', 'Alkalinity of ash', 'Magnesium', 'Total phenols', 'Flavanoids', 'Nonflavanoid phenols', 'Proanthocyanins', 'Color intensity', 'Hue', 'OD280/OD315 of diluted wines', 'Proline'

Class weights:

$$\pi_1 = 43/130 = 0.33, \pi_2 = 51/130 = 0.39, \pi_3 = 36/130 = 0.28$$

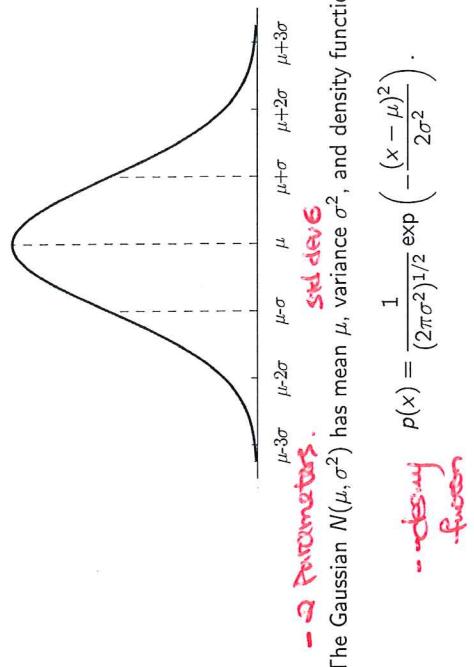
Need distributions P_1, P_2, P_3 , one per class. **Fitting a distribution to each class.**

Base these on a single feature: 'Alcohol'.

(B) Features: 13-dimensional dataset, so we need \rightarrow a distribution over 13-dimensional space

Alcohol \downarrow reduces to 1-dimensional space; fit one distribution to the alcohol level from basket 1 (so 13 (QQ)).

Q: What one-dimensional distribution should we use?
The univariate Gaussian (default).



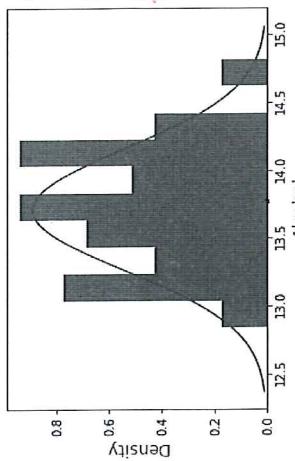
$$p(x) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

The Gaussian $N(\mu, \sigma^2)$ has mean μ , variance σ^2 , and density function

The distribution for winery 1

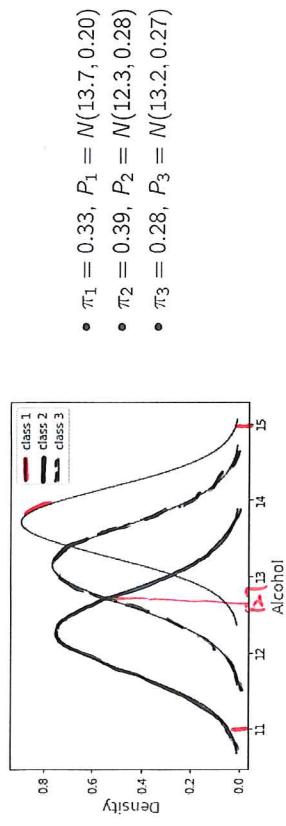
Single feature: 'Alcohol'

the lens is the density that we fit for winery number 1.



Mean $\mu = 13.72$, Standard deviation $\sigma = 0.44$ (variance 0.20)

All three wineries



To classify x : Pick the j with highest $\pi_j P_j(x)$

$\pi_1 = 0.33$ ~~highest~~

Test error: $14/48 = 29\%$

$\lambda = 12.7$ ~~(W2 > W3)~~

$$z = \frac{x - \mu}{\sigma}$$

$$z_1 = \frac{12.7 - 13.7}{0.2} = -5$$

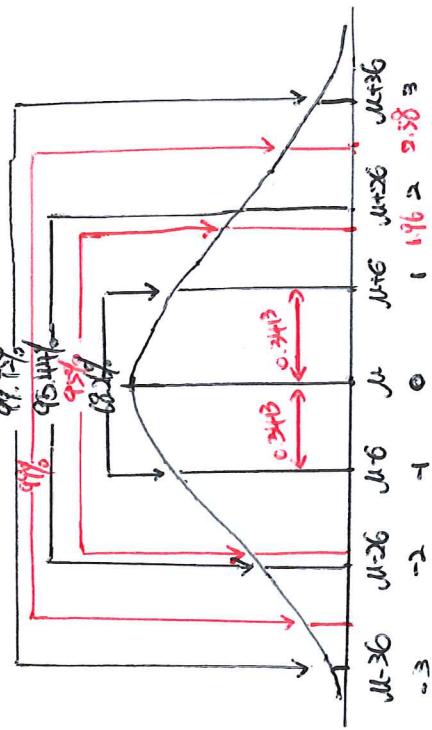
z for the z-table find $0.3413 (34.13\%)$

$$34.13\% \times 0.28 = 68.26\%$$

The percentage of people
that score between 1 and -1
is 68.26% of the population.

- 68.26% - in total, observations within plus or minus one standard deviation from the mean ($\pm \frac{\sigma}{2}$) ; ~~percentile class between~~
- 95.44% - ~~percentile class between~~
- 99.72% - ~~percentile class between~~

Normal Distribution



- 0.5% - corresponds to 95% of the total distribution;
- 0.05% left on the z-table (0.05).
- 0.5% - corresponds to 99% of the total distribution.



Random variables

Weeks 4
15/Jan/2018.

Probability review II:

Random variables, expectation, and variance

Roll two dice. Let X be their sum.

$$\begin{aligned}\text{outcome} = (1, 1) &\Rightarrow X = 2 \\ \text{outcome} = (1, 2) \text{ or } (2, 1) &\Rightarrow X = 3\end{aligned}$$

Probability space:

- Sample space: $\Omega = \{1, 2, 3, 4, 5, 6\} \times \{1, 2, 3, 4, 5, 6\}$.
- Each outcome equally likely.

Random variable X lies in $\{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$.

A random variable (r.v.) is a function defined on a probability space.

It is a mapping from Ω (outcomes) to \mathbb{R} (numbers).
We'll use capital letters for r.v.'s.

the information in the underlying probability space.
Q: How we can determine the distribution of a random variable

Topics we'll cover

Roll a die.
Define $X = 1$ if die is ≥ 3 , otherwise $X = 0$.

① What is a random variable?

② Expected value

③ Variance and standard deviation

$$X \in \{0, 1\}$$

$$\Pr(X=0) = \Pr(\text{die} = (1, 2)) = \frac{2}{6} = \frac{1}{3}$$

$$\Pr(X=1) = \Pr(\text{die} = 3, 4, 5, 6) = \frac{4}{6} = \frac{2}{3}$$

Expected value, or mean

Expected value of a random variable X : The expected x is simply the value of x if we repeated the experiment

$$\mathbb{E}(X) = \sum_x x \cdot \Pr(X=x)$$

Roll a die. Let X be the number observed. This is a weighted average of all possible values that X can take.

What is $\mathbb{E}(X)$?
 $X \in \{1, 2, 3, 4, 5, 6\}$

$$\begin{aligned}\mathbb{E}(X) &= 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + 3 \cdot \frac{1}{6} + 4 \cdot \frac{1}{6} + 5 \cdot \frac{1}{6} + 6 \cdot \frac{1}{6} \\ &= 3.5\end{aligned}$$

A property of expected values

How is the average of a set of numbers affected if:
• You double the numbers?
• You increase each number by 1?
• The average also increase by 1?

Summary: Let X be any random variable.
If $V = aX + b$ (any constants a, b), then $\mathbb{E}(V) = a\mathbb{E}(X) + b$

- A new random variable V .
- V is some constant times X plus some other constant.

Linearity property

Another example

A biased coin has heads probability p . Let X be 1 if heads, 0 if tails. What is $\mathbb{E}(X)$?

$$\begin{aligned}\mathbb{E}(X) &= 0 \cdot \Pr(X=0) + 1 \cdot \Pr(X=1) \\ &= 0 \cdot (1-p) + 1 \cdot p \\ &= p.\end{aligned}$$

Variance

Can summarize an r.v. X by its mean, μ . But this doesn't capture the spread of X :

- Two distributions with same mean μ .
- Left, tightly centered around its mean;
- Right, more dispersed.
- How can we capture expected value of abs distance from μ ?
- Variance: $\text{var}(X) = \mathbb{E}((X - \mu)^2)$, where $\mu = \mathbb{E}(X)$
- Standard deviation $\sqrt{\text{var}(X)}$: Average square dist from mean.
- Roughly, the average amount by which X differs from its mean.



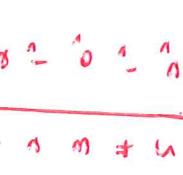
Variance: example

Choose X uniformly at random from $\{1, 2, 3, 4, 5\}$.

$$\mathbb{E}(x) = 1 \cdot \frac{1}{5} + 2 \cdot \frac{1}{5} + 3 \cdot \frac{1}{5} + 4 \cdot \frac{1}{5} + 5 \cdot \frac{1}{5} = 3 = \text{mean.}$$

$$\text{var}(x) = \mathbb{E}[(x - \mu)^2] = (0)^2 \cdot \frac{1}{5} + (1)^2 \cdot \frac{1}{5} + (2)^2 \cdot \frac{1}{5} + (3)^2 \cdot \frac{1}{5} + (4)^2 \cdot \frac{1}{5} = 2.$$

$$\begin{array}{c} x \\ \times \\ (x - \mu)^2 \\ \hline 1 & 2^2 \\ 2 & 1^2 \\ 3 & 0^2 \\ 4 & 1^2 \\ 5 & 2^2 \end{array}$$



$$\mathbb{E}(X^2) = 1^2 \cdot \frac{1}{5} + 2^2 \cdot \frac{1}{5} + 3^2 \cdot \frac{1}{5} + 4^2 \cdot \frac{1}{5} + 5^2 \cdot \frac{1}{5} = 11$$

$$\text{Var}(x) = \mathbb{E}(X^2) - \mu^2 = 11 - 3^2 = 2.$$

Alternative formula for variance

Variance: $\text{var}(X) = \mathbb{E}((X - \mu)^2)$, where $\mu = \mathbb{E}(X)$

Another way to write it: $\text{var}(X) = \mathbb{E}(X^2) - \mu^2$ minus mean squared.

Example: Choose X uniformly at random from $\{1, 2, 3, 4, 5\}$.

- expected value of x squared.

$$\mathbb{E}(X^2) = 1^2 \cdot \frac{1}{5} + 2^2 \cdot \frac{1}{5} + 3^2 \cdot \frac{1}{5} + 4^2 \cdot \frac{1}{5} + 5^2 \cdot \frac{1}{5} = 11$$

$$\text{Var}(x) = \mathbb{E}(X^2) - \mu^2 = 11 - 3^2 = 2.$$

Variance: properties

Variance: $\text{var}(X) = \mathbb{E}((X - \mu)^2)$, where $\mu = \mathbb{E}(X)$

- Variance is always ≥ 0 non-negative number.
- How is the variance affected if:
 - You increase each number by 1?
 - You double each number? Variance is doubled. Variance is $\times 4$ but the spread remains the same).
- Summary: If $V = aX + b$ then $\text{var}(V) = a^2 \text{var}(X)$

Week 5

Probability review III: Measuring dependence

Independent random variables

Random variables X, Y are **independent** if $\Pr(X = x, Y = y) = \Pr(X = x)\Pr(Y = y)$.

• X being equal to x has ~~either~~ either
made it more likely or less likely
that Y being equals to y .

Pick a card out of a standard deck.
 $X =$ suit and $Y =$ number.

$$\Pr(X = \heartsuit, Y = 5) = \frac{1}{52}$$

$$\Pr(X = \heartsuit) = \frac{1}{4}$$

$$\Pr(Y = 5) = \frac{1}{13}$$

- when dealing with multiple random variables

e.g., a patient's heart rate, blood pressure ... the easiest situation is all variables are independent.

- Then don't worry about the interactions between them.

Topics we'll cover

Random variables X, Y are **independent** if $\Pr(X = x, Y = y) = \Pr(X = x)\Pr(Y = y)$.

Flip a fair coin 10 times.
 $X =$ # heads and $Y =$ last toss.

- Q. Are X and Y independent?
- ① When are two random variables independent?
 - ② Qualitatively assessing dependence
 - ③ Quantifying dependence: covariance and correlation

$$\Pr(X = 10 \cdot Y = T) = 0$$

$$\Pr(X = 10) = \left(\frac{1}{2}\right)^{10}$$

$$\Pr(Y = T) = \frac{1}{2}$$

∴ Not independent variables.

Independent random variables

Random variables X, Y are independent if $\Pr(X = x, Y = y) = \Pr(X = x)\Pr(Y = y)$.

$X, Y \in \{-1, 0, 1\}$, with these probabilities:

	X	Y	P_{XY}
X	-1	-1	0.4
	0	0	0.16
	1	1	0.03

To check if those are independent, we have to check every entry in this table of joint distribution.

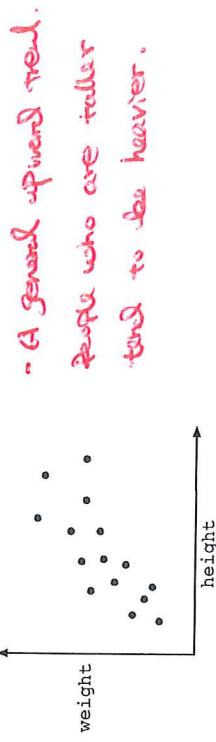
$$\Pr(X = -1 \text{ and } Y = 1) = 0.03 \\ = 0.8 \times 0.3$$

\therefore independent X and Y .

Positive correlation

	X	Y	P_{XY}
X	-1	-1	0.5
	0	0	0.2
	1	1	0.3

H, W are positively correlated



- A general upward trend.
- people who are taller tend to be heavier.

This also implies $\mathbb{E}[HW] > \mathbb{E}[H]\mathbb{E}[W]$.

The average H times W .

Average H times Average W .

Dependence

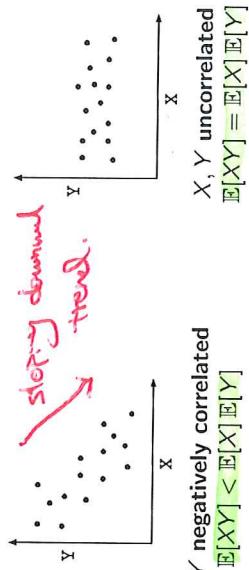
Example: Pick a person at random, and take

$$\begin{aligned} H &= \text{height} \\ W &= \text{weight} \end{aligned}$$

Independence would mean

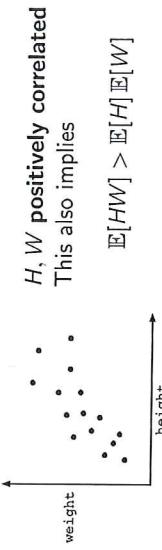
$$\Pr(H = h, W = w) = \Pr(H = h)\Pr(W = w).$$

Not accurate: height and weight will be positively correlated.



X, Y negatively correlated
 $\mathbb{E}[XY] < \mathbb{E}[X]\mathbb{E}[Y]$

Types of correlation



H, W positively correlated
This also implies

$$\mathbb{E}[HW] > \mathbb{E}[H]\mathbb{E}[W]$$

X, Y uncorrelated
 $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$

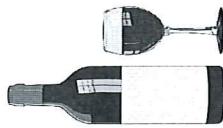
- we would expect height and weight to be positively correlated.



The winery prediction problem

Which winery is it from, 1, 2, or 3?

week 0-8 91 Two classes
Two-dimensional generative modeling with the bivariate Gaussian



one feature used

Using one feature ('Alcohol'), error rate is 29%.

What if we use two features? - error rate drops

- look at Gaussian in 2 dimensions. Allow us to model the dependence between features.

Topics we'll cover

Training set obtained from 130 bottles

- Winery 1: 43 bottles
 - Winery 2: 51 bottles
 - Winery 3: 36 bottles
 - For each bottle, 13 features:
 - 'Alcohol', 'Malic acid', 'Ash', 'Alcalinity of ash', 'Magnesium', 'Total phenols', 'Flavanoids', 'Nonflavanoid phenols', 'Proanthocyanins', 'Color intensity', 'Hue', 'OD280/OD315 of diluted wines', 'Proline'
- Also, a separate test set of 48 labeled points.

This time: 'Alcohol' and 'Flavanoids'.

Q. why it might be helpful to throw in a 3rd feature?

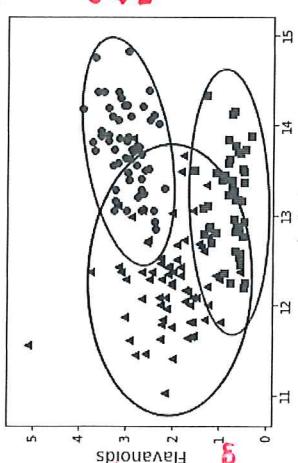
- Fit a Gaussian to alcohol level from class 1;
- " " class 2;
- " " class 3.

Alcohol

Why it helps to add features

Better separation between the classes!

- The 3 Gaussians aren't separated very well; one feature alone is not a great basis for drug classification.



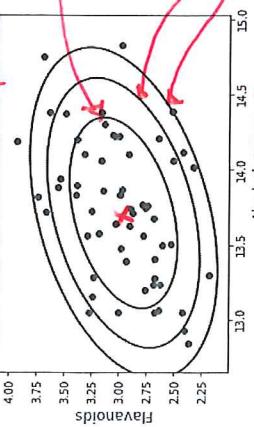
- Two features may be better for class fusion.

Error rate drops from 29% to 8%.

- Error rate drops after ~ 3%.

The bivariate Gaussian

- 443 two-dimensional points from class 1.



- Fit a Gaussian distribution with 2 parameters.

Model class 1 by a bivariate Gaussian, parametrized by:

$$\text{mean } \mu = \begin{pmatrix} 13.7 \\ 3.0 \end{pmatrix} \text{ and covariance matrix } \Sigma = \begin{pmatrix} 0.20 & 0.06 \\ 0.06 & 0.12 \end{pmatrix}$$

- Center of distribution.
- Highest density.

- As move away from mean, the density drops off, and it drops off with the ellipsoidal contours.

The bivariate (2-d) Gaussian

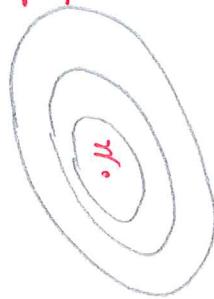
Suppose X_1 has mean μ_1 and X_2 has mean μ_2 .

Can measure dependence between them by their covariance:

- $\text{cov}(X_1, X_2) = \mathbb{E}[(X_1 - \mu_1)(X_2 - \mu_2)] = \mathbb{E}[X_1 X_2] - \mu_1 \mu_2$
- Maximized when $X_1 = X_2$, in which case it is $\text{var}(X_1)$.
- It is at most $\text{std}(X_1)\text{std}(X_2)$.
- Cor = negative \rightarrow negative correlation
- Cor = 0 \rightarrow non correlation

$$\left\{ \begin{array}{l} \Sigma_{11} = \text{var}(X_1) \\ \Sigma_{22} = \text{var}(X_2) \\ \Sigma_{12} = \Sigma_{21} = \text{cov}(X_1, X_2) \end{array} \right\}$$

- μ : the highest density
- the shape of ellipsoids is given by sigma matrix



- Density is highest at the mean, falls off in ellipsoidal contours.

Center of a Gaussian

$$\text{mean } \mu = \begin{pmatrix} 13.7 \\ 3.0 \end{pmatrix}$$

- variance along the 1st direction
- variance along the 2nd direction

