

Schedule

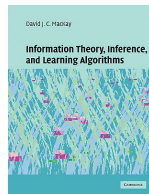
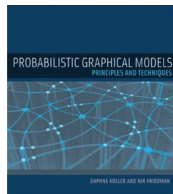
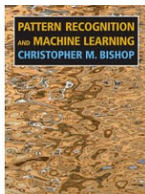
	Refresher of Probabilities
	Introduction to Probabilistic Graphical Models
	Probabilistic Inference
	Learning Conditional Random Fields
	MAP Prediction / Energy Minimization
	Learning Structured Support Vector Machines

Links to slide download: `http://pub.ist.ac.at/~chl/courses/PGM_W16/`

Password for ZIP files (if any): `pgm2016`

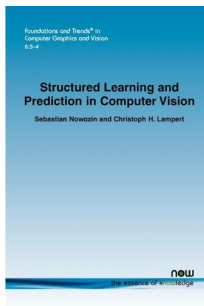
Email for questions, suggestions or typos that you found: `chl@ist.ac.at`

For the curious ones...



- ▶ Bishop, **Pattern Recognition and Machine Learning**, Springer New York, 2006, ISBN-13: 978-0387310732
- ▶ Koller, Friedman, **Probabilistic Graphical Models: Principles and Techniques**, The MIT Press, 2009, ISBN-13: 978-0262013192
- ▶ MacKay, **Information Theory, Inference and Learning Algorithms**, Cambridge University Press, 2003, ISBN-13: 978-0521642989

Older tutorial...



Parts published in

- ▶ Sebastian Nowozin, Chrsitoph H. Lampert, "**Structured Learning and Prediction in Computer Vision**", Foundations and Trends in Computer Graphics and Vision, now publisher, <http://www.nowpublishers.com/>
- ▶ available as PDF on my homepage

Success Stories of Machine Learning



Robotics



Time Series Prediction

amazon.com

Recommended for You

Body Science
Our Price: **\$9.99**
Used & new from \$9.99
[See all buying options](#)

Because you purchased...

The Black Swan: Second Edition: The Impact of the Highly Improbable: With a

Social Networks



Language Processing

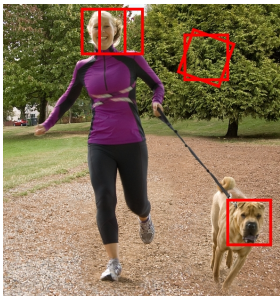
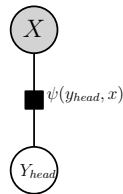


Healthcare



Natural Sciences

All of these require dealing with Structured Data

Human Pose \mathcal{Y} Image $x \in \mathcal{X}$ Example y_{head} 

Head detector

- Idea: Have a head classifier (CNN, SVM, ...) $\psi(y_{head}, x) \in \mathbb{R}_+$

Refresher: Probabilities

Random Variables

A **random variable** is a variable that randomly takes one of its possible values:

- ▶ the number of photons reaching a CCD chip
- ▶ the text of the next email I will receive
- ▶ the position of an atom in a molecule

Some notation: we will write

- ▶ random variables with capital letters, *e.g.* X
- ▶ the set of possible values it can take with curly letters, *e.g.* \mathcal{X}
- ▶ any individual value it can take with lowercase letters, *e.g.* x

How likely each value $x \in \mathcal{X}$ is specified by a *probability distribution*.

There are, slightly different, possibilities:

- ▶ \mathcal{X} is discrete (typically finite),
- ▶ \mathcal{X} is continuous.

Discrete Random Variables

For discrete \mathcal{X} (e.g. $\mathcal{X} = \{0, 1\}$):

- ▶ $p(X = x)$ is the probability that X takes the value $x \in \mathcal{X}$.
If it's clear which variable we mean, we'll just write $p(x)$.
- ▶ for example, rolling a die, $p(X = 3) = p(3) = 1/6$
- ▶ we write $x \sim p(x)$ to indicate that the distribution of X is $p(x)$

For things to make sense, we need

$$0 \leq p(x) \leq 1 \quad \text{for all } x \in \mathcal{X} \quad \text{(positivity)}$$

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

Example: English words

- ▶ X_{word} : pick a word randomly from an English text. Is it "word"?
- ▶ $\mathcal{X}_{\text{word}} = \{\text{true}, \text{false}\}$

$$p(X_{\text{the}} = \text{true}) = 0.05$$

$$p(X_{\text{the}} = \text{false}) = 0.95$$

$$p(X_{\text{horse}} = \text{true}) = 0.004$$

$$p(X_{\text{horse}} = \text{false}) = 0.996$$

Continuous Random Variables

For continuous \mathcal{X} (e.g. $\mathcal{X} = \mathbb{R}$):

- ▶ probability that X takes a value in the set M is

$$\Pr(X \in A) = \int_M p(x) dx$$

- ▶ we call $p(x)$ the **probability density over x**

For things to make sense, we need:

$$p(x) \geq 0 \quad \text{for all } x \in \mathcal{X} \quad \text{(positivity)}$$

$$\int_{\mathcal{X}} p(x) = 1 \quad \text{(normalization)}$$

Note: for convenience of notation, we use the notation of discrete random variable everywhere.

Joint probabilities

Probabilities can be assigned to more than one random variable at a time:

- ▶ $p(X = x, Y = y)$ is the probability that $X = x$ and $Y = y$ (at the same time)

joint probability

Example: English words

Pick three consecutive English words: $X_{word}, Y_{word}, Z_{word}$

- ▶ $p(X_{the} = \text{true}, Y_{horse} = \text{true}) = 0.00080$
- ▶ $p(X_{horse} = \text{true}, Y_{the} = \text{true}) = 0.00001$
- ▶ $p(X_{probabilistic} = \text{true}, Y_{graphical} = \text{true}, Y_{model} = \text{true}) = 0.000000045$

Conditional probabilities

One random variable can contain information about another one:

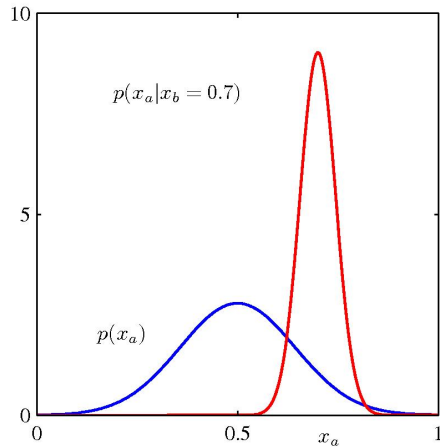
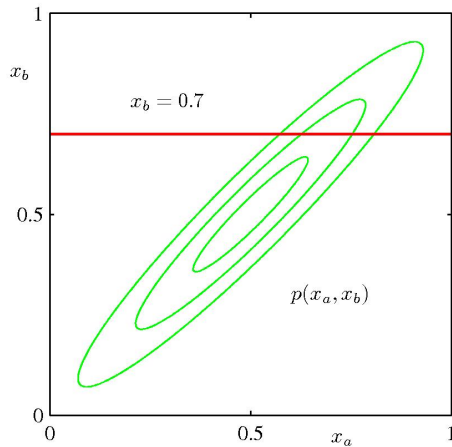
- ▶ $p(X = x | Y = y)$: **conditional probability**
what is the probability of $X = x$, if we already know that $Y = y$?
- ▶ $p(X = x)$: **marginal probability**
what is the probability of $X = x$, without any additional information?
- ▶ conditional probabilities can be computed from joint and marginal:

$$p(X = x | Y = y) = \frac{p(X = x, Y = y)}{p(Y = y)} \quad (\text{not defined if } p(Y = y) = 0)$$

Example: English text

- ▶ $p(X_{the} = \text{true}) = 0.05$
- ▶ $p(X_{the} = \text{true} | Y_{horse} = \text{true}) = \frac{0.0008}{0.004} = 0.20$
- ▶ $p(X_{the} = \text{true} | Y_{the} = \text{true}) = \frac{0.0003}{0.05} = 0.006$

Illustration



joint (level sets), marginal, conditional probability

Bayes rule (Bayes theorem)

Bayes rule

Most famous formula in probability:
$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}$$

Bayes rule (Bayes theorem)

Bayes rule

Most famous formula in probability: $p(A|B) = \frac{p(B|A)p(A)}{p(B)}$

Formally, nothing spectacular: direct consequence of definition of conditional probability.

$$p(A|B) = \frac{p(A, B)}{p(B)} = \frac{p(B|A)p(A)}{p(B)}$$

Bayes rule (Bayes theorem)

Bayes rule

Most famous formula in probability: $p(A|B) = \frac{p(B|A)p(A)}{p(B)}$

Formally, nothing spectacular: direct consequence of definition of conditional probability.

$$p(A|B) = \frac{p(A, B)}{p(B)} = \frac{p(B|A)p(A)}{p(B)}$$

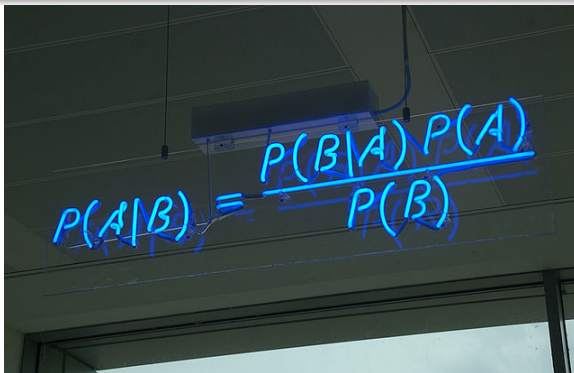
Nevertheless very useful at least for two situations:

- ▶ when A and B have different role, so $p(A|B)$ is intuitive but $p(B|A)$ is not
 - ▶ $A = \text{age}$, $B = \{\text{smoker, nonsmoker}\}$
 - $p(A|B)$ is the age distribution amongst smokers and nonsmokers
 - $p(B|A)$ is the probability that a person of a certain age smokes
- ▶ the information in B help us to update our knowledge about A : $p(A) \mapsto p(A|B)$

Bayes rule (Bayes theorem)

Bayes rule

Most famous formula in probability:
$$p(A|B) = \frac{p(B|A)p(A)}{p(B)}$$



Dependence/Independence

Not every random variable is informative about every other.

- ▶ We say X is independent of Y (write: $X \perp\!\!\!\perp Y$) if

$$p(X = x, Y = y) = p(X = x)p(Y = y) \quad \text{for all } x \in \mathcal{X} \text{ and } y \in \mathcal{Y}$$

- ▶ equivalent (if defined):

$$p(X = x|Y = y) = p(X = x), \quad p(Y = y|X = x) = p(Y = y)$$

Dependence/Independence

Not every random variable is informative about every other.

- ▶ We say **X is independent of Y (write: $X \perp\!\!\!\perp Y$)** if

$$p(X = x, Y = y) = p(X = x)p(Y = y) \quad \text{for all } x \in \mathcal{X} \text{ and } y \in \mathcal{Y}$$

- ▶ equivalent (if defined):

$$p(X = x|Y = y) = p(X = x), \quad p(Y = y|X = x) = p(Y = y)$$

Other random variables can influence the independence:

- ▶ **X and Y are conditionally independent given Z (write $X \perp\!\!\!\perp Y|Z$)** if

$$p(X = x, Y = y|Z = z) = p(X = x|Z = z)p(Y = y|Z = z)$$

- ▶ equivalent (if defined):

$$p(x|y, z) = p(x|z), \quad p(y|x, z) = p(y|z)$$

Example: rolling dice

Let X and Y be the outcome of independently rolling two dice and let $Z = X + Y$ be their sum.

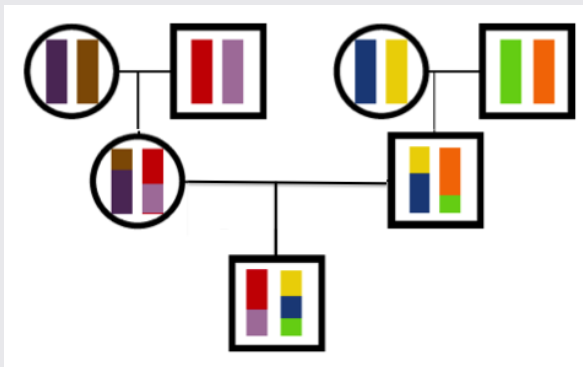
- ▶ X and Y are independent
- ▶ X and Z are not independent, Y and Z are not independent
- ▶ conditioned on Z , X and Y are not independent anymore
(for fixed $Z = z$, X and Y can only take certain value combinations)

Example: toddlers

Let X be the height of a toddler, Y the number of words in its vocabulary and Z its age.

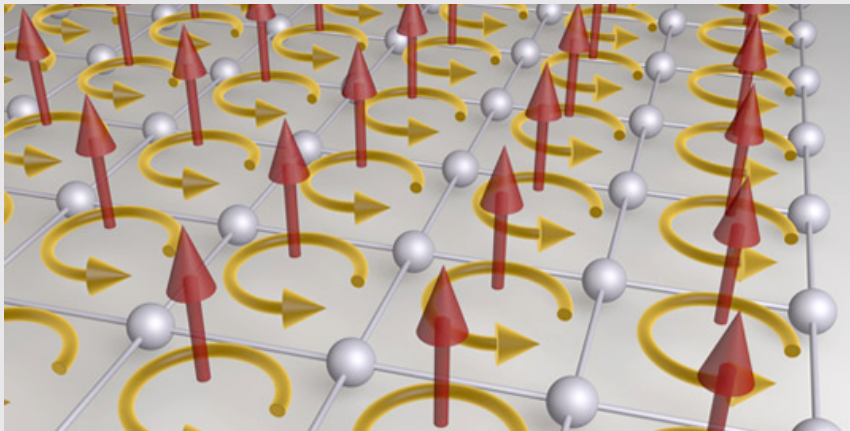
- ▶ X and Y are not independent: overall, toddlers who are taller know more words
- ▶ however, X and Y are conditionally independent given Z :
at a fixed age, toddlers' growth and vocabulary develop independently

Example



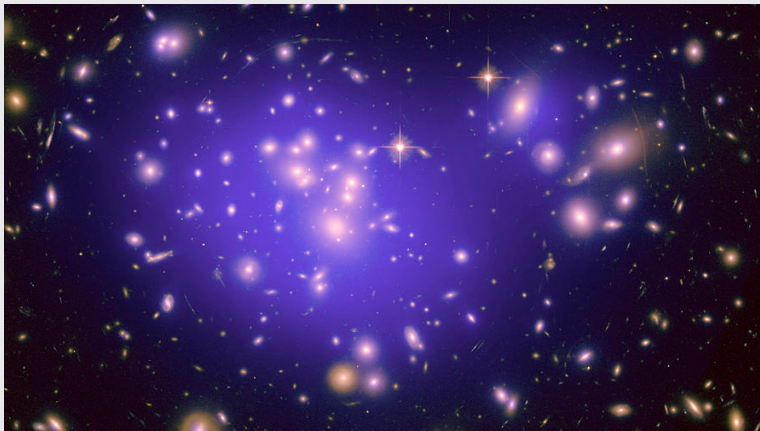
- ▶ X = your genome
- ▶ Y_1, Y_2 = your parents' genomes
- ▶ Z_1, Z_2, Z_3, Z_4 = your grandparents' genomes

Discrete Random Fields



Magnetic spin in each atoms of a crystal: $X_{i,j}$ for $i, j \in \mathbb{Z}$

Continuous Random Fields



Distribution of matter in the universe: X_p for $p \in \mathbb{R}^3$

Expected value

We apply a function to (the values of) one or more random variables:

$$\blacktriangleright f(x) = x^2 \quad \text{or} \quad f(x_1, x_2, \dots, x_k) = \frac{x_1 + x_2 + \dots + x_k}{k}$$

The **expected value** or **expectation** of a function f with respect to a probability distribution is the weighted average of the possible values:

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

In short, we just write $\mathbb{E}_x[f(x)]$ or $\mathbb{E}[f(x)]$ or $\mathbb{E}[f]$ or $\mathbb{E}f$.

Example: rolling dice

Let X be the outcome of rolling a die and let $f(x) = x$

$$\mathbb{E}_{x \sim p(x)}[f(x)] = \mathbb{E}_{x \sim p(x)}[x] = \frac{1}{6}1 + \frac{1}{6}2 + \frac{1}{6}3 + \frac{1}{6}4 + \frac{1}{6}5 + \frac{1}{6}6 = 3.5$$

Expected value

Example: rolling dice

X_1, X_2 : the outcome of rolling two dice independently, $f(x_1, x_2) = x_1 + x_2$

$$\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] =$$

Expected value

Example: rolling dice

X_1, X_2 : the outcome of rolling two dice independently, $f(x_1, x_2) = x_1 + x_2$

$$\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] =$$

Straight-forward computation: 36 options for (x_1, x_2) , each has probability $\frac{1}{36}$

Expected value

Example: rolling dice

X_1, X_2 : the outcome of rolling two dice independently, $f(x_1, x_2) = x_1 + x_2$

$$\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] =$$

Straight-forward computation: 36 options for (x_1, x_2) , each has probability $\frac{1}{36}$

$$\begin{aligned}\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] &= \sum_{x_1, x_2} p(x_1, x_2)(x_1 + x_2) \\ &= \frac{1}{36}(1 + 1) + \frac{1}{36}(1 + 2) + \frac{1}{36}(1 + 3) + \dots \\ &\quad + \frac{1}{36}(2 + 1) + \frac{1}{36}(2 + 2) + \dots + \frac{1}{36}(6 + 6) \\ &= \frac{252}{36} = 7\end{aligned}$$

Expected value

Example: rolling dice

X_1, X_2 : the outcome of rolling two dice independently, $f(x_1, x_2) = x_1 + x_2$

$$\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] = 7$$

Straight-forward computation: 36 options for (x_1, x_2) , each has probability $\frac{1}{36}$

$$\begin{aligned}\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] &= \sum_{x_1, x_2} p(x_1, x_2)(x_1 + x_2) \\ &= \frac{1}{36}(1 + 1) + \frac{1}{36}(1 + 2) + \frac{1}{36}(1 + 3) + \dots \\ &\quad + \frac{1}{36}(2 + 1) + \frac{1}{36}(2 + 2) + \dots + \frac{1}{36}(6 + 6) \\ &= \frac{252}{36} = 7\end{aligned}$$

Expected value

Example: rolling dice

X_1, X_2 : the outcome of rolling two dice independently, $f(x_1, x_2) = x_1 + x_2$

$$\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] =$$

Sometimes a good strategy: count how often each value occurs and sum over values

$s = (x_1 + x_2)$	2	3	4	5	6	7	8	9	10	11	12
count n_s	1	2	3	4	5	6	5	4	3	2	1

Expected value

Example: rolling dice

X_1, X_2 : the outcome of rolling two dice independently, $f(x_1, x_2) = x_1 + x_2$

$$\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] = 7$$

Sometimes a good strategy: count how often each value occurs and sum over values

$s = (x_1 + x_2)$	2	3	4	5	6	7	8	9	10	11	12
count n_s	1	2	3	4	5	6	5	4	3	2	1

$$\begin{aligned}\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] &= \sum_{x_1, x_2} p(x_1, x_2)(x_1 + x_2) = \sum_s \frac{n_s}{n} s \\ &= \frac{1}{36} 2 + \frac{2}{36} 3 + \frac{3}{36} 4 + \frac{4}{36} 5 + \dots + \frac{2}{36} 11 + \frac{1}{36} 12 = \frac{252}{36} = 7\end{aligned}$$

Properties of expected values

Example: rolling dice

X_1, X_2 : the outcome of rolling two dice independently, $f(x_1, x_2) = x_1 + x_2$

$$\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] =$$

Properties of expected values

Example: rolling dice

X_1, X_2 : the outcome of rolling two dice independently, $f(x_1, x_2) = x_1 + x_2$

$$\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] =$$

The expected value has a useful property: it is *linear* in its argument.

- ▶ $\mathbb{E}_{x \sim p(x)}[f(x) + g(x)] = \mathbb{E}_{x \sim p(x)}[f(x)] + \mathbb{E}_{x \sim p(x)}[g(x)]$
- ▶ $\mathbb{E}_{x \sim p(x)}[\lambda f(x)] = \lambda \mathbb{E}_{x \sim p(x)}[f(x)]$

If a random variables does not show up in a function, we can ignore the expectation operation with respect to it

- ▶ $\mathbb{E}_{(x, y) \sim p(x, y)}[f(x)] = \mathbb{E}_{x \sim p(x)}[f(x)]$

Properties of expected values

Example: rolling dice

X_1, X_2 : the outcome of rolling two dice independently, $f(x_1, x_2) = x_1 + x_2$

$$\begin{aligned}\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] &= \mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[x_1 + x_2] \\ &= \mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[x_1] + \mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[x_2] \\ &= \mathbb{E}_{x_1 \sim p(x_1)}[x_1] + \mathbb{E}_{x_2 \sim p(x_2)}[x_2] = 3.5 + 3.5 = \mathbf{7}\end{aligned}$$

The expected value has a useful property: it is *linear* in its argument.

- ▶ $\mathbb{E}_{x \sim p(x)}[f(x) + g(x)] = \mathbb{E}_{x \sim p(x)}[f(x)] + \mathbb{E}_{x \sim p(x)}[g(x)]$
- ▶ $\mathbb{E}_{x \sim p(x)}[\lambda f(x)] = \lambda \mathbb{E}_{x \sim p(x)}[f(x)]$

If a random variables does not show up in a function, we can ignore the expectation operation with respect to it

- ▶ $\mathbb{E}_{(x, y) \sim p(x, y)}[f(x)] = \mathbb{E}_{x \sim p(x)}[f(x)]$

Example: rolling dice

- ▶ we roll one die
- ▶ X_1 : number facing up, X_2 : number facing down
- ▶ $f(x_1, x_2) = x_1 + x_2$

$$\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] =$$

Example: rolling dice

- ▶ we roll one die
- ▶ X_1 : number facing up, X_2 : number facing down
- ▶ $f(x_1, x_2) = x_1 + x_2$

$$\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] = 7$$

Answer 1: explicit calculation with dependent X_1 and X_2

$$p(x_1, x_2) = \begin{cases} \frac{1}{6} & \text{for combinations } (1, 6), (2, 5), (3, 4), (4, 3), (5, 2), (6, 1) \\ 0 & \text{for all other combinations.} \end{cases}$$

$$\begin{aligned} \mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] &= \sum_{(x_1, x_2)} p(x_1, x_2)(x_1 + x_2) \\ &= 0(1 + 1) + 0(1 + 2) + \dots + \frac{1}{6}(1 + 6) + 0(2 + 1) + \dots = 6 \cdot \frac{7}{6} = 7 \end{aligned}$$

Example: rolling dice

- ▶ we roll one die
- ▶ X_1 : number facing up, X_2 : number facing down
- ▶ $f(x_1, x_2) = x_1 + x_2$

$$\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] = 7$$

Answer 2: use properties of expectation as earlier

$$\begin{aligned}\mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[f(x_1, x_2)] &= \mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[x_1 + x_2] \\ &= \mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[x_1] + \mathbb{E}_{(x_1, x_2) \sim p(x_1, x_2)}[x_2] \\ &= \mathbb{E}_{x_1 \sim p(x_1)}[x_1] + \mathbb{E}_{x_2 \sim p(x_2)}[x_2] = 3.5 + 3.5 = 7\end{aligned}$$

The rules of probability take care of dependence, etc.

Some expected values show up so often that they have special names.

Variance

The **variance** of a random variable X is the expected squared deviation from its mean

$$\text{Var}(X) = \mathbb{E}_x[(x - \mathbb{E}_x[x])^2]$$

also

$$\text{Var}(X) = \mathbb{E}_x[x^2] - (\mathbb{E}_x[x])^2 \quad (\text{exercise})$$

The variance

- ▶ measures how much the random variable *fluctuates* around its mean
- ▶ is invariant under addition

$$\text{Var}(X + a) = \text{Var}(X) \quad \text{for } a \in \mathbb{R}.$$

- ▶ scales with the square of multiplicative factors

$$\text{Var}(\lambda X) = \lambda^2 \text{Var}(X) \quad \text{for } \lambda \in \mathbb{R}.$$

More intuitive:

Standard deviation

The **standard deviation** of a random variable is the square root of the its variance.

$$\text{Std}(X) = \sqrt{\text{Var}(X)}$$

The standard deviation

- ▶ is invariant under addition

$$\text{Std}(X + a) = \text{Std}(X) \quad \text{for } a \in \mathbb{R}.$$

- ▶ scales with the absolute value of multiplicative factors

$$\text{Std}(\lambda X) = |\lambda| \text{Std}(X) \quad \text{for } \lambda \in \mathbb{R}.$$

For two random variables at a time, we can test if their fluctuations around the mean are consistent or not

Covariance

The **covariance** of two random variables X and Y is the expected value of the product of their deviations from their means

$$\text{Cov}(X, Y) = \mathbb{E}_{(x,y) \sim p(x,y)}[(x - \mathbb{E}_x[x])(y - \mathbb{E}_y[y])]$$

The covariance

- ▶ of a random variable with itself it its variance, $\text{Cov}(X, X) = \text{Var}(X)$
- ▶ is invariant under addition: $\text{Cov}(X + a, Y) = \text{Cov}(X, Y) = \text{Cov}(X, Y + a)$ for $a \in \mathbb{R}$.
- ▶ scales linearly under multiplications:
 $\text{Cov}(\lambda X, Y) = \lambda \text{Cov}(X, Y) = \text{Cov}(X, \lambda Y)$ for $\lambda \in \mathbb{R}$.
- ▶ is 0, if X and Y are independent, but can be 0 even if for dependent X and Y (exercise)

If we do not care about the scales of X and Y , we can normalize by their standard deviations:

Correlation

The **correlation coefficient** of two random variables X and Y is their covariance divided by their standard deviations

$$\text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\text{Std}(X) \text{Std}(Y)} = \frac{\mathbb{E}_{(x,y) \sim p(x,y)}[(x - \mathbb{E}_x[x])(y - \mathbb{E}_y[y])]}{\sqrt{\mathbb{E}_x(x - \mathbb{E}_x[x])^2} \sqrt{\mathbb{E}_y(y - \mathbb{E}_y[y])^2}}$$

The correlation

- ▶ always has values in the interval $[-1, 1]$
- ▶ is invariant under addition: $\text{Cov}(X + a, Y) = \text{Cov}(X, Y) = \text{Cov}(X, Y + a)$ for $a \in \mathbb{R}$.
- ▶ is invariant under multiplication with positive constants
 $\text{Corr}(\lambda X, Y) = \text{Corr}(X, Y) = \text{Corr}(X, \lambda Y)$ for $\lambda > 0$.
- ▶ inverts its sign under multiplication with negative constants
 $\text{Corr}(-\lambda X, Y) = -\text{Corr}(X, Y) = \text{Corr}(X, -\lambda Y)$ for $\lambda > 0$.
- ▶ is 0, if X and Y are independent, but can be 0 even if for dependent X and Y (exercise)

