

Probability Concepts

- We will make use of the following concepts.
 - Basic probability in discrete spaces, events
 - Joint probability
 - Conditional probability
 - Independence (and conditional independence)
 - Marginal probability (marginalization)
 - Probability in continuous spaces (probability density functions)
- To learn/review them see supplementary chapter in the book posted on the web site or your favorite web text or video resource

Probabilistic Fitting

- Generative probabilistic model
 - Tells a story about how stochastic data comes to be
 - Darts fall around the center of the board, but where exactly?
 - Consider a model with parameters, θ
 - Consider an observation, x_i
 - We denote the probability of seeing x_i under the model by:

$$p(x_i | \Theta)$$

↑
Read “given” or “conditioned on”
Restricts to the case of θ

Defined by $P(A|B) = \frac{P(A,B)}{P(B)}$

Probabilistic Fitting

- Multiple observations
 - Suppose we have multiple observations, in a vector \mathbf{x}
 - What is the probability of \mathbf{x} ?
- If observations are independent then probability is the product of the individual observations
 - Essentially a definition, but it is consistent with intuition
 - The observations are conditionally independent **given** the model
- So, the probability of \mathbf{x} is then:

$$p(\mathbf{x} | \Theta) = \prod p(x_i | \Theta)$$

Probabilistic Fitting

- So, given the model, we have the probability of observing the data

$$p(\mathbf{x} | \Theta) = \prod p(x_i | \Theta)$$

- But what we really want is the probability of the model (parameters) given the data!
- Bayes rule comes to the rescue!

Bayes Rule

- Bayes rule:
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$
- Proof
$$P(A,B) = P(B|A)P(A) = P(A|B)P(B)$$
- With our notation:
$$P(\Theta|\mathbf{x}) = \frac{P(\mathbf{x}|\Theta)P(\Theta)}{P(\mathbf{x})}$$

likelihood function
for the parameters

prior probability (often
taken to be uniform)

$$P(\Theta|\mathbf{x}) = \frac{P(\mathbf{x}|\Theta)P(\Theta)}{P(\mathbf{x})}$$

posterior probability

normalizer, often is
not of interest

Common special case
 $P(\Theta|\mathbf{x}) \propto P(\mathbf{x}|\Theta)$

Know the words in **red**

Probabilistic Fitting

- If we assume **uniform** prior, then we can find the posterior density for the parameters by:

$$p(\Theta|\mathbf{x}) \propto p(\mathbf{x}|\Theta)$$

- Now the objective is to find the parameters Θ such that this *likelihood* is maximum
- Note--this is the same as finding the parameters which minimize the **negative log likelihood**

Probabilistic fitting with independence and uniform prior

Finding the “best” model under simple circumstances

maximize $p(\Theta|\mathbf{x})$ (one definition of best Θ)

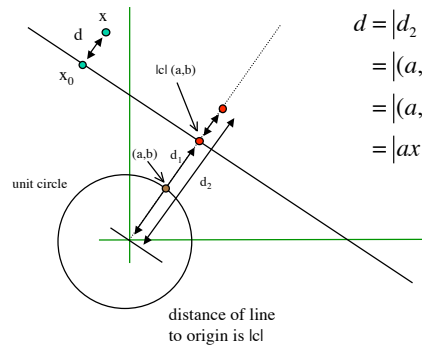
maximize $p(\mathbf{x}|\Theta)$ (by Bayes rule, uniform prior)

minimize $-\log(p(\mathbf{x}|\Theta))$ (log is monotonic increasing)

minimize $-\log\left(\prod p(x_i|\Theta)\right)$ (by independence)

minimize $-\sum \log(p(x_i|\Theta))$ (high school math)

- Back to lines: $ax+by+c=0$ where $a^2+b^2=1$
- Distance squared from (x,y) to this line is $(ax+by+c)^2$



$$\begin{aligned} d &= |d_2 - d_1| \\ &= |(a,b) \cdot x - (a,b) \cdot x_0| \\ &= |(a,b) \cdot x + c| \\ &= |ax + by + c| \end{aligned}$$

- **Generative model** for lines: Choose point on line, and then, with probability proportional to $p(d)$, **normally distributed** (Gaussian), go a distance d from the line.

- Now the probability of an observed (x,y) is given by

$$p((x,y) | \Theta) \propto \exp\left(-\frac{(ax + by + c)^2}{2\sigma^2}\right)$$

Lines

Convenient formula for line
 $ax+by+c=0$
where $a^2+b^2=1$

(x_D, y_D)

$$d^2 = (ax_D + by_D + c)^2$$

This is the generative model
It tells us $P(\text{data} | \text{model})$

$$p((x_D, y_D) | \Theta) \propto \exp\left(-\frac{(ax_D + by_D + c)^2}{2\sigma^2}\right)$$

We have the probability density of the observed (x,y) given by

$$p((x,y) | \Theta) \propto \exp\left(-\frac{(ax + by + c)^2}{2\sigma^2}\right)$$

The negative log is

$$\frac{(ax + by + c)^2}{2\sigma^2}$$

And the negative log likelihood of multiple observations is

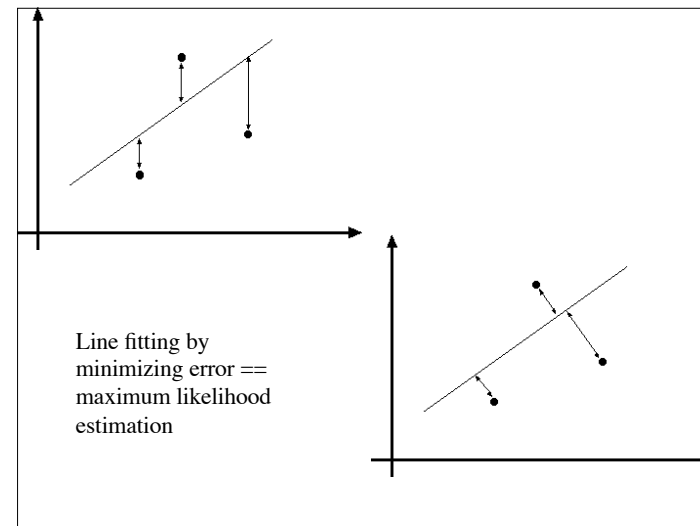
$$\frac{1}{2\sigma^2} \sum_i (ax_i + by_i + c)^2$$

From the previous slide, we had that the negative log likelihood of multiple observations is given by

$$-\frac{1}{2\sigma^2} \sum_i (ax_i + by_i + c)^2 \quad (\text{where } a^2 + b^2 = 1)$$

This should be recognizable as homogeneous least squares

Thus we have shown that least squares is maximum likelihood estimation under normality (Gaussian) error statistics!



Fitting curves other than lines

- In principle, an easy generalization
 - For Gaussian error statistics, Euclidean distance is a good measure
 - The probability of obtaining a point, given a curve, is given by a negative exponential of distance squared
- In practice, this can be hard
 - It can be difficult to compute the distance between a point and a curve
 - Circles, ellipses, and a few others are not too hard
 - Otherwise, craft an approximation
 - §15.3 has more