

## Probabilistic Interpretation

### - Probability Refresh

- **Sample space S:** set of possible outcomes
- **Random Event E:** subset of S
- **Random Variable:** maps event E to a real value (**denoted by Capitals**)
- **Conditional Probability**
  - **Events:**  $P(E|F) = \frac{P(E \cap F)}{P(F)}$  when  $P(F) > 0$
  - **RVs:**  $P(X = x|Y = y) = \frac{P(X=x \text{ and } Y=y)}{P(Y=y)}$
- **Chain Rule:**  $P(X = x \text{ and } Y = y) = P(X = x|Y = y)P(Y = y)$
- **Consequences of the Chain Rule**
  - **Marginalisation:** Suppose RV Y takes values in  $\{y_1, y_2, \dots, y_n\}$  Then
    - $P(X = x) = P(X = x \text{ and } Y = y_1) + \dots + P(X = x \text{ and } Y = y_n)$
    - $= \sum_{i=1}^n P(X = x|Y = y_i)P(Y = y_i)$
  - **Bayes Rule:**  $P(X = x|Y = y) = \frac{P(Y=y|X=x)P(X=x)}{P(Y=y)}$
  - **Independence:** Random variables X and Y are independent if...
    - $P(X = x \text{ and } Y = y) = P(X = x)P(Y = y)$
    - For all x and y, in which case:  $P(X = x|Y = y) = P(X = x)$
- **Continuous-valued random variables**
  - $P(X = x) = 0$  for continuous-valued random variables, so we consider intervals instead:  $P(a \leq X \leq b)$
  - $F_Y(y) := P(Y \leq y)$  is the **cumulative distribution function (CDF)** and  $P(a < Y \leq b) = F_Y(b) - F_Y(a)$
  - For a continuous-valued random variable, Y, there exists a **probability density function**  $f_Y(y) \geq 0$  such that:
    - $F_Y(y) = \int_{-\infty}^y f_Y(t)dt$
    - And so...
    - $P(a < Y \leq b) = \int_{-\infty}^b f_Y(t)dt - \int_{-\infty}^a f_Y(t)dt = \int_a^b f_Y(t)dt$
    - **The probability density function f(y) for random variable Y is not a probability (it can take values greater than 1) - the area under the PDF is the probability  $P(a < Y \leq b)$**
  - $\int_{-\infty}^{\infty} f(y)dy = 1$  (since  $\int_{-\infty}^{\infty} f(y)dy = F_Y(\infty) = P(Y \leq \infty) = 1$ )
- **CDF for X and Y:**  $F_{XY}(x, y) = P(X \leq x \text{ and } Y \leq y)$ 
  - Well defined for both continuous and discrete valued RVs

- When X and Y are continuous-valued RVs there exists a PDF

$$f_{XY}(x,y) \geq 0 \text{ such that: } F_{XY}(x,y) = \int_{-\infty}^x \int_{-\infty}^y f_{XY}(u,v) du dv$$

- Define conditional PDF:  $f_{X|Y}(x|y) = \frac{f_{XY}(x,y)}{f_Y(y)}$

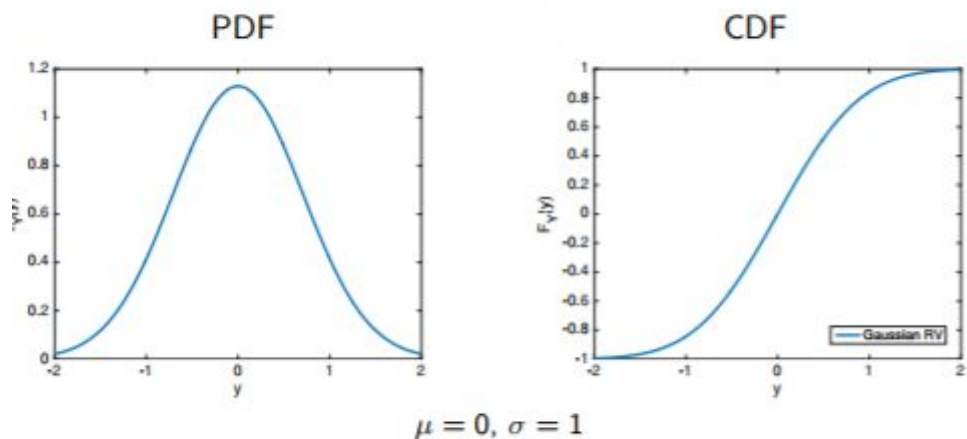
- Then the chain rule holds for PDFs:

$$f_{XY}(x,y) = f_{X|Y}(x|y)f_Y(y) = f_{Y|X}(y|x)f_X(x)$$

- **So marginalisation, Bayes rule and independence carry over to PDFs similarly to discrete-valued RVs**

- Y is a Normal or Gaussian RV  $Y \sim N(\mu, \sigma^2)$  when it has the PDF:

$$f_Y(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$



- $E[Y] = \mu, \text{Var}(Y) = \sigma^2$
- Symmetric about  $\mu$  and defined for all real-valued  $x$

## - Probabilistic Interpretation of Linear Regression

- Assume output Y is generated by:

$$Y = \theta^T x + M = h_\theta(x) + M$$

- Where  $h_\theta(x) = \theta^T x$  and M is Gaussian noise with mean 0 and variance 1

- So training data d is:

$$\{(x^{(1)}, h_\theta(x^{(1)}) + M^{(1)}), (x^{(2)}, h_\theta(x^{(2)}) + M^{(2)}), \dots, (x^{(m)}, h_\theta(x^{(m)}) + M^{(m)})\}$$

- Where  $M^{(1)}, M^{(2)}, \dots, M^{(m)}$  are independent RVs each of which is Gaussian with mean 0 and variance 1

- A Gaussian RV Z with mean  $\mu$  and variance  $\sigma^2$  has PDF:

$$f_Z(z) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(z-\mu)^2}{2\sigma^2}}$$

- So we are assuming:

$$f_M(m) = \frac{1}{\sqrt{2\pi}} e^{-\frac{m^2}{2}}, f_Y(y) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(y-h_\theta(x))^2}{2}}$$

- The **likelihood**  $f_{D|\Theta}(d|\theta)$  of the training data  $d$  is therefore:

$$f_{D|\Theta}(d|\theta) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}} e^{-\frac{(y^{(i)} - h_{\theta}(x^{(i)}))^2}{2}}$$

- Taking logs:  $\log f_{D|\Theta}(d|\theta) = \log \frac{1}{\sqrt{2\pi}} - \sum_{i=1}^m \frac{(y^{(i)} - h_{\theta}(x^{(i)}))^2}{2}$

- And the maximum likelihood estimate of  $\theta$  maximises:

$$\max_{\theta} - \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)}))^2$$

$$\text{I.e. it minimises: } \min_{\theta} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)}))^2$$

### - Why do we care about probabilistic interpretation?

- Probability is the reasoning about uncertainty, it would be strange if machine learning algorithms didn't make sense from a probability perspective
- Casting an ML approach within a statistical framework clarifies the assumptions we (implicitly) make e.g. in Linear Regression:
  - **Noise is additive:**  $Y = \theta^T x + M$
  - **Noise on each observation is independent and identically distributed**
  - **Noise is Gaussian** - this is what drives our usage of square loss. Changing the noise model would lead to a different loss function
- Allows us to utilise the results and approaches of probability/statistics, and perhaps gain new insights. E.g. in linear regression:
  - Without regularisation, our estimate of  $\theta$  is the maximum likelihood estimate. Would a MAP (Maximum A Posteriori) estimate be more/less useful?

### - Probabilistic Interpretation of Logistic Regression

- Assume:

$$P(Y = y|\theta, x) = \frac{1}{1 + e^{-y\theta^T x}}$$

- And recall  $y = 1$  or  $y = -1$  only

- The **likelihood** of training data  $d$  is:  $f_{D|\Theta}(d|\theta) = \prod_{i=1}^m \frac{1}{1 + e^{-y\theta^T x}}$

- **Taking logs:**  $\log f_{D|\Theta}(d|\theta) = \sum_{i=1}^m \log \frac{1}{1 + e^{-y\theta^T x}}$

- **And the maximum likelihood estimate of  $\theta$  minimises:**

$$-\sum_{i=1}^m \log \frac{1}{1 + e^{-y\theta^T x}} = \sum_{i=1}^m \log(1 + e^{-y\theta^T x})$$

- Since  $-\log(z) = \log(1/z)$
- **The probabilistic formulation of logistic regression provides us with new insight:**
  - $P(Y = y|\theta, x) = \frac{1}{1+e^{-y\theta^T x}}$
- So in addition to our prediction  $h_\theta(x) = \text{sign}(\theta^T x)$  we also have a confidence in the prediction:  $\frac{1}{1+e^{-y\theta^T x}}$ 
  - When  $\frac{1}{1+e^{-y\theta^T x}}$  is close to 1 we are confident, close to zero we are less confident
- **Probabilistic Interpretation of Regularisation**

$$P(\Theta = \vec{\theta} | D = d) = \frac{P(D = d | \Theta = \vec{\theta}) P(\Theta = \vec{\theta})}{P(D = d)}$$

posterior                  likelihood                  prior

- **Bayes Rule:**
  - **Likelihood:** probability of seeing the data  $d$ , given the model with parameter  $\Theta = \theta$  where  $\theta$  is a vector
  - **Prior:** Before seeing any data what is our belief of the model... i.e. **what is probability of parameter values  $\Theta$**
  - **Posterior:** after seeing the data, what is our belief about probability of parameter values  $\Theta$  now that we have seen the data
- **Maximum A Posteriori (MAP):** estimate of vector  $\theta$  is value that maximises  $P(\Theta = \theta | D = d)$
- **Maximum Likelihood estimation:** Select value that maximises  $P(D = d | \Theta = \theta)$
- **Taking logs in Bayes rule:**
  - $\log P(\Theta = \theta | D = d) = \log P(D = d | \Theta = \theta) + \log P(\Theta = \theta) - \log P(D = d)$
  - Can drop the  $\log P(D = d)$  as  $d$  is fixed, so we select  $\theta$  to maximise:
    - $\log P(D = d | \Theta = \theta) + \log P(\Theta = \theta)$
  - Or for continuous-valued RVs:
    - $\log f_{D|\Theta}(D = d | \Theta = \theta) + \log f_\Theta(\Theta = \theta)$
- **Ridge regression variant of linear regression:**

- $Y = \Theta x + M$ ,  $M \sim N(0, 1)$  as before.
- $\Theta_j, \sim N(0, \sigma^2)$  (this is our prior on  $\theta_j$ ),  $j = 1, \dots, n$
- log-likelihood:  $-\sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2$
- log-prior:  $-\theta_j^2/\sigma^2$
- So MAP estimate selects  $\theta$  to maximise:

$$-\sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 - \sum_{j=1}^n \theta_j^2/\sigma^2$$

i.e. to minimise:

$$\sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 + \sum_{j=1}^n \theta_j^2/\sigma^2$$