Actual Bayesian computat°:

Let's now introduce a prior $P(\theta)$, i.e. a guess about our idea of what is $\theta$.

Recall: MLE: $\text{argmax}_\theta \left( P(X \mid \theta) \right) = \theta^*_{MLE}$

likelihood of the data $\quad\hookrightarrow$ leads to the usual $\hat{\mu} = \frac{1}{N} \sum_i^N x_i$, etc

Now MAP: Ma<u>x</u>imum a <u>Posteriori</u>:

$$\theta^*_{MAP} = \text{argmax}_\theta \; P(\theta \mid X) \qquad \text{likelihood of the parameters}$$

$$= \text{argmax}_\theta \left( \frac{\overbrace{P(X \mid \theta)}^{\text{likelihood of data}} \; P(\theta)}{P(X)} \right) \to \text{prior}$$

$$\hookrightarrow \text{evidence (un\underline{reach}able)}$$

$$= \text{argmax}_\theta \left( \underbrace{\log P(X \mid \theta)}_{\text{as before}} + \underbrace{\log P(\theta)}_{\text{prior}} \right)$$

Example (simple): $\quad X \sim (x_i)_{1 \ldots N}, \quad X \sim \mathcal{N}(\mu, \sigma)$

N examples. prior: $\mu \sim \mathcal{N}(0, \tau)$.

$$\hookrightarrow P(\theta) = \frac{1}{\sqrt{2\pi}\,\tau} \cdot e^{-\frac{1}{2} \cdot \frac{\mu^2}{\tau^2}} \quad \theta$$

$$\theta^*_{MAP} = \text{argmax}_\theta \left( \log \prod_{i=1}^N \left( \frac{1}{\sqrt{2\pi}\,\sigma} e^{-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}} \right) + \log \left( \frac{1}{\sqrt{2\pi}\,\tau} e^{-\frac{1}{2} \frac{\mu^2}{\tau^2}} \right) \right)$$

$$\frac{\delta \ell}{\delta \mu} = \frac{\partial}{\partial \mu} \sum_{i=1}^N \left[ \log\left( \frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{1}{2} \frac{(x-\mu)^2}{\sigma^2} \right] + cte - \frac{1}{2} \frac{\mu^2}{\tau^2}$$

$$= 0 - \frac{2}{2} \sum_i^N \frac{x_i - \mu}{\sigma^2} - \frac{2}{2} \frac{\mu}{\tau^2}$$

$$\Longrightarrow \sum_{i=1}^N \frac{x_i}{\sigma^2} = N \frac{\mu}{\sigma^2} - \frac{\mu}{\tau^2} \Longleftrightarrow \mu \left( 1 - \frac{\sigma^2}{N\tau^2} \right) = \frac{1}{N} \sum_i^N x_i$$

So we find: $\qquad\qquad\qquad \hat{\sigma}_{MAP} = \hat{\sigma}_{MLE}$ (here)

$$\hat{\mu}_{MAP} = \hat{\mu}_{MLE} \times \frac{1}{1 - \frac{\sigma^2}{N\tau^2}} \qquad \text{(this is just 1 example)}$$

Remarks : $\tau = \infty \Rightarrow$ flat prior $\Rightarrow$ no change

$\qquad\qquad N \to \infty \Rightarrow \infty$ data $\Rightarrow$ prior becomes
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ irrelevant.

$\qquad\qquad \sigma \to 0 \Rightarrow$ data little spread $\Rightarrow$ prior irrelevant

We see that the prior acts as a regularizat°
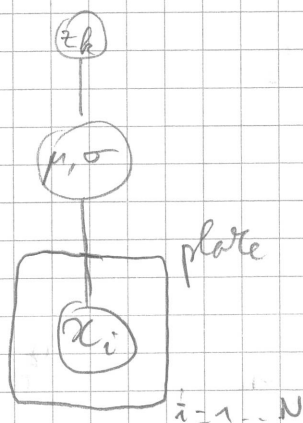
Many interpretations of regul° as Bayesian priors. ☺

[ When Sparse distrib° (eg words), need a prior to
avoid singularities

[ N$_{train}$ small $\to$ try to use your knowledge
$\qquad\qquad\qquad\qquad\qquad\qquad$ (put it in the prior)

Words distro : Many words $\to$ Not all appear
in the train set (corpus) : $\hat{\mu}_{word \# ...} = 0$

But, we expect they can appear (eg test set) :

So, need a prior with $p\left(word = "\left[\substack{some\ infrequent \\ word}\right]"\right) > 0$



prior on the prior : assume $\tau \sim \mathcal{N}(m, B)$