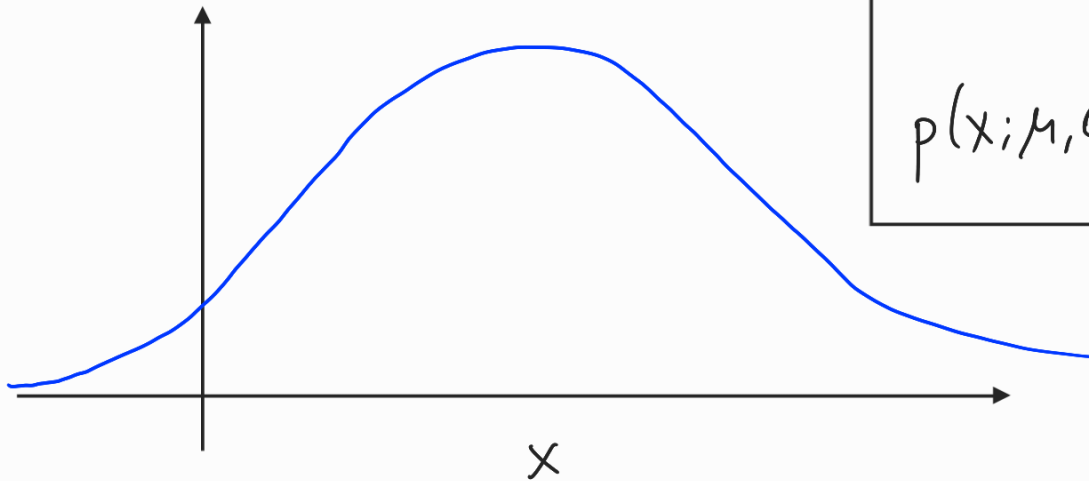


Gaussian distribution:

$x \in \mathbb{R}$ with mean = μ and variance = σ^2 ;

$$x \sim \mathcal{N}(\mu, \sigma^2)$$



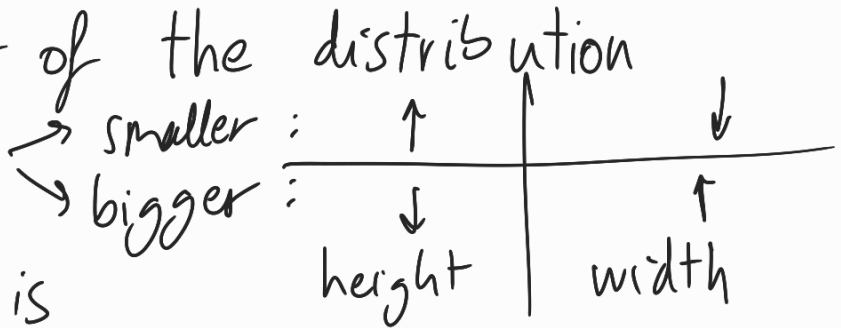
$$p(x; \mu, \sigma^2) = \frac{e^{(-\frac{(x-\mu)^2}{2\sigma^2})}}{\sqrt{2\pi} \sigma}$$

$\mu \Rightarrow$ sets the center of the distribution

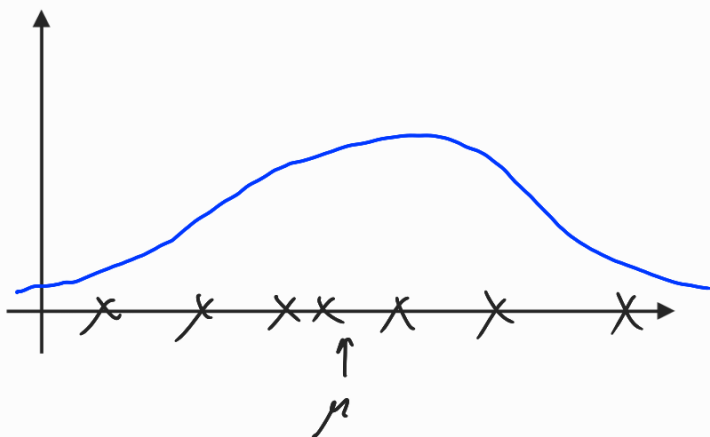
$\sigma^2 \Rightarrow$ height and width

That's like that,
because the AREA is

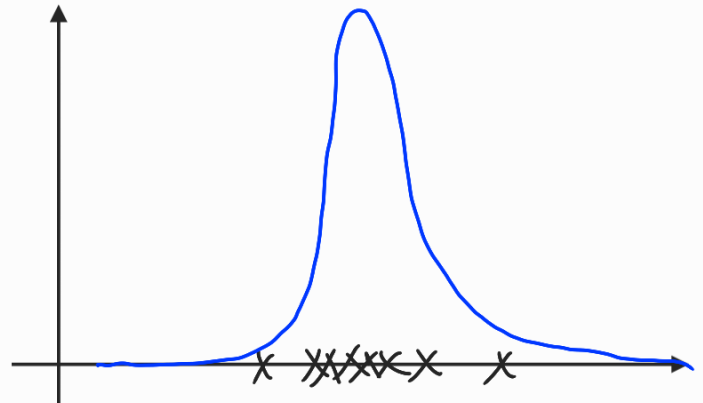
ALWAYS the same.



As the σ^2 grows:



As it decreases:



Real values are more far from μ .

Density estimation:

Training set: $\{x^{(1)}, \dots, x^{(m)}\} \forall x \in \mathbb{R}^n$

With its corresponding $\{\mu_1, \sigma_1^2\}, \dots, \{\mu_m, \sigma_m^2\}$

$$p(x) = \prod_{i=1}^m p(x_i; \mu_i, \sigma_i^2) = \prod_{i=1}^m \frac{e^{-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}}}{\sqrt{2\pi}\sigma_i}$$

Anomaly detection algorithm:

1. Choose features that may be anomalous

2. Fit parameters μ and σ

3. Given new x , compute $p(x) < \epsilon$

Example:

1000 good

20 anomalous

Training set:

6000
≠

Cross validation:

2000
≠

Test set:

2000

Anomalous
0

10
≠

10

How to evaluate model?

- Accuracy: NO, very skewed

- F1 score: YES

Choose ϵ : with cross validation, iteratively maximize F_1 .

	Normal			Anomalous	Many different types of anomalies
Anomaly detection	↑	↑	↑	↑	
Supervised learning	↓	↓	↓	↓	

Multivariate Gaussian distribution: used when we need more complex functions than just \odot .

Being able then to create ellipses with Σ' and moving them with μ ; $\forall \mu \in \mathbb{R}^n$, $\Sigma' \in \mathbb{R}^{n \times n}$

1. Fit model $p(x)$

$$\hookrightarrow \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad \hookrightarrow \Sigma' = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

2. Given new x , compute:

$$p(x) = \frac{e^{-\frac{1}{2} (x - \mu)^T \Sigma'^{-1} (x - \mu)}}{(2\pi)^{n/2} |\Sigma'|^{1/2}}$$

$p(x) < \epsilon$: Anomaly

Original model	Multivariate model
<ul style="list-style-type: none">- Manually create features- Cheaper- OK if training size small (m small)	<ul style="list-style-type: none">- Automatically captures corrs.- Expensive- Number samples $>$ number features (if not, $\nexists \Sigma'^{-1}$)