National Technical University of Ukraine
"Igor Sikorsky Kyiv Polytechnic Institute"
Institute of Physics and Technology

# Lecture 4
# Generative Models for Real-World Systems. Frequentist Concept

Dmytro Progonov,

PhD, Associate Professor,

Department Of Physics and Information Security Systems

# Content

- Frequentist statistics concept;

- Estimator risk minimization procedures;

- Estimator's properties;

- Pathologies of frequentist statistics.

Generative Models for Real-World Systems.
Frequentist Concept

2/11

# Frequentist statistics concept

**_Frequentist_** (**_classical_** or **_orthodox_**) **_statistics_** are based on the concept of a **_sampling distribution_**. This is the distribution that an **_estimator_** $\delta$ has when applied to the multiple dataset $\{\mathcal{D}_1, \mathcal{D}_2 \cdots \mathcal{D}_K\}$ sampled from the true but unknown distributions.

**_Parameter estimate_** $- \hat{\theta} = \delta(\mathcal{D})$.

In classical (frequentist) decision theory, there is a loss function and a likelihood, but there is no prior and hence no posterior or posterior expected loss. Thus *there is no automatic way to deriving an optimal estimator, unlike the Bayesian case*.

Having chosen an estimator $\delta$, we define its expected **_loss_** or **_risk_** as

$$R(\theta^*, \delta) \triangleq \mathbb{E}_{p(\widetilde{\mathcal{D}}|\theta^*)}\big[L\big(\theta^*, \delta(\widetilde{\mathcal{D}})\big)\big] = \int L\big(\theta^*, \delta(\widetilde{\mathcal{D}})\big) p(\widetilde{\mathcal{D}}|\theta^*) d\,\widetilde{\mathcal{D}},$$

where $\widetilde{\mathcal{D}}$ is a data sampled from "nature's distribution", which is represented by parameter $\theta^*$.

Generative Models for Real-World Systems.
Frequentist Concept

3/11

# Estimator risk minimization procedures. Bayes risk

One approach to convert risk $R(\theta^*, \delta)$ into a single measure of quantity, $R(\delta)$, which does not depend on knowing $\theta^*$, is to put a prior on $\theta^*$, and then to define ***Bayes (integrated) risk***:

$$R_B(\delta) \triangleq \mathbb{E}_{p(\theta^*)}[R(\theta^*, \delta)] = \int R(\theta^*, \delta) p(\theta^*) d\theta^*$$

A ***Bayes estimator*** or ***Bayes decision rule*** is one which minimizes the expected risk:

$$\delta_B \triangleq \underset{\delta}{\operatorname{argmin}} R_B(\delta)$$

*Theorem*: A Bayes estimator can be obtained by minimizing the posterior expected loss for each **x**.

*Theorem (Wald, 1950)*: Every admissible decision rule is a Bayes decision rule with resoect to some, possibly improper, prior distribution.

Generative Models for Real-World Systems.
Frequentist Concept

4/11

# Estimator risk minimization procedures. Minimax risk

An alternative approach to Bayes estimator, that does not require choice of prior, is usage of minimax rule:

$$\delta_{MM} \triangleq \underset{\delta}{\text{argmin}}\, R_{max}(\delta),$$

where $R_{max}(\delta) \triangleq \underset{\delta}{\max}\, R(\theta^*, \delta)$ — the maximum risk of an estimator.

Minimax estimator is very pessimistic — one can show that all minimax estimators are equivalent to Bayes estimators under a ***least favorable prior***.

Generative Models for Real-World Systems.
Frequentist Concept

5/11

# Estimator risk minimization procedures. Empirical risk minimization

We define the empirical risk as follows:

$$R_{emp}(\mathcal{D}, \delta) \triangleq R\big(p_{emp}(\cdot \,|\mathcal{D}), \delta\big) = \frac{1}{N} \sum_{i=1}^{N} L(y_i, \delta(\mathbf{x}_i))$$

A _**empirical risk minimization**_ is task of finding a decision procedure to minimize the empirical risk:

$$\delta_{ERM} \triangleq \underset{\delta}{\operatorname{argmin}} \, R_{emp}(\mathcal{D}, \delta)$$

Possible approaches for solving the mentioned task:

_**Regularized risk minimization**_:

$$\acute{R}(\mathcal{D}, \delta) = R_{emp}(\mathcal{D}, \delta) + \lambda C(\delta)$$

_**Structural risk minimization**_:

$$\begin{cases} \acute{\delta}_\lambda = \underset{\delta}{\operatorname{argmin}} \big[R_{emp}(\mathcal{D}, \delta) + \lambda C(\delta)\big] \\ \acute{\lambda} = \underset{\lambda}{\operatorname{argmin}} \, \acute{R}(\acute{\delta}_\lambda) \end{cases}$$

Generative Models for Real-World Systems.
Frequentist Concept

6/11

# Estimator's properties

_**Consistent estimator**_ – if it eventually recovers the true parameters that generated the data as the sample size goes to infinity, i.e. $\hat{\theta}(\mathcal{D}) \to \theta^*$ as $|\mathcal{D}| \to \infty$ (convergence in probability sense).

_**Unbiased estimator**_ provides the null bias of parameters measurement:

$$bias\left(\hat{\theta}(\cdot)\right) = \mathbb{E}_{p(\mathcal{D}|\theta_*)}\big[\hat{\theta}(\mathcal{D}) - \theta_*\big] \to 0$$

_**Minimum variance estimator**_ achieves the lower bound on the variance of any unbiased estimator:

_Cramer-Rao inequality_: Let $X_1, X_2 \cdots X_n \sim p(X|\theta_0)$ and $\hat{\theta} = \hat{\theta}(x_1, x_2 \cdots x_n)$ be an unbiased estimator of $\theta_0$. Then, under various smoothness assumption on $p(X|\theta_0)$ we have

$$var[\hat{\theta}] \geq \frac{1}{nI(\theta_0)},$$

where $I(\theta_0)$ is Fisher information matrix.

$$MSE = variance + bias^2$$

Generative Models for Real-World Systems.
Frequentist Concept

7/11

# Pathologies of frequentist statistics.
# Counter-intuitive behavior of confidence intervals

A frequentist confidence interval for some parameter $\theta$ is defined y expression:

$$C_\alpha(\theta) = (l, u): P\big(l(\widetilde{\mathcal{D}}) \leq \theta \leq u(\widetilde{\mathcal{D}})|\widetilde{\mathcal{D}}{\sim}\theta\big) = 1 - \alpha.$$

_In Bayesian statistics_, we condition on what is known – namely the observed data, $\mathcal{D}$ – and average on what is unknown, namely the parameter $\theta$.

_In frequentist statistics_, we do exactly the opposite: we condition on what is unknown – namely the true parameter value $\theta$ – and average over hypothetical future dataset $\widetilde{\mathcal{D}}$.

Generative Models for Real-World Systems.
Frequentist Concept

8/11

# Pathologies of frequentist statistics.
# P-values considered harmful

Suppose we want to decide whether to accept/reject some baseline model (*null hypothesis*). In frequentist statistics, it is standard to compute the **_p-value_**, which is defined as probability (under the null) of observing some test statistics $f(\mathcal{D})$:

$$pvalue(\mathcal{D}) \triangleq P\big(f(\widetilde{\mathcal{D}}) \geq f(\mathcal{D})|\widetilde{\mathcal{D}} \sim H_0\big)$$

Given the p-value, we reject the null hypothesis iff the p-value is less than some threshold (such as $\alpha = 0.05$). If we reject it, we say the difference between observed and expected test statistics is **_statistically significant_** at level $\alpha$.

This procedure guarantees that our expected type I (false positive) error rate is at most $\alpha$. This sometimes is interpreted as saying that _frequentist hypothesis testing is very conservative_ (unlikely to accidently reject the null hypothesis). In fact the opposite is the case – because this method only worries about trying to reject the null, it can never gather evidence in favor of the null, no matter how large the sample size.

Another problem with *p-values* is that their _computation depends on decision you make about when to stop collecting data_, even if these decisions do not change the data you actually observed.

Generative Models for Real-World Systems.
Frequentist Concept

9/11

# Pathologies of frequentist statistics.
# The likelihood principle

The fundamental reason for mentioned pathologies is that frequentist statistics violates:

***Likelihood principle*** – interference should be based on likelihood of the observed data, not based on hypothetical future data that you have not observed.

***Sufficiency principle*** – sufficient statistics contains all the relevant information about an unknown parameter.

***Weak conditionality*** – interference should be based on events that happened, not which might have happened.

Generative Models for Real-World Systems.
Frequentist Concept

10/11

# Conclusion

- Concept of frequentist statistics was presented;

- Estimators and theirs properties were considered;

- Procedures for estimator risk minimization were shown;

- Limitations of practical usage the frequentist statistics were considered.

Generative Models for Real-World Systems. Frequentist Concept

11/11