

FAIRNESS RISK MEASURES

FAIRNESS RISK MEASURES



Robert C. Williamson



Australian
National
University

Aditya Menon



Google AI

LOSS FUNCTIONS – OUTCOME CONTINGENT UTILITIES





LOSS FUNCTIONS - OUTCOME CONTINGENT UTILITIES

- ▶ Wald's abstraction: a loss function

$$\ell : \underset{\substack{\text{Label} \\ \text{space}}}{Y} \times \underset{\substack{\text{Action} \\ \text{space}}}{A} \rightarrow \mathbb{R}_+ \cup \{+\infty\} =: \overline{\mathbb{R}}$$



LOSS FUNCTIONS - OUTCOME CONTINGENT UTILITIES

- ▶ Wald's abstraction: a loss function

$$\ell : \underset{\substack{\text{Label} \\ \text{space}}}{Y} \times \underset{\substack{\text{Action} \\ \text{space}}}{A} \rightarrow \mathbb{R}_+ \cup \{+\infty\} =: \overline{\mathbb{R}}$$

- ▶ $a \mapsto \ell(y, a)$ is an outcome contingent utility



LOSS FUNCTIONS - OUTCOME CONTINGENT UTILITIES

- ▶ Wald's abstraction: a loss function

$$\ell : \underset{\substack{\text{Label} \\ \text{space}}}{Y} \times \underset{\substack{\text{Action} \\ \text{space}}}{A} \rightarrow \mathbb{R}_+ \cup \{+\infty\} =: \overline{\mathbb{R}}$$

- ▶ $a \mapsto \ell(y, a)$ is an outcome contingent utility
- ▶ Learning goal: **expected** risk minimisation

$$\min_{f \in \mathcal{F}} \mathbb{E}_{(X, Y) \sim P} \ell(Y, f(X))$$



LOSS FUNCTIONS – OUTCOME CONTINGENT UTILITIES

- ▶ Wald's abstraction: a loss function

$$\ell : \underset{\substack{\text{Label} \\ \text{space}}}{Y} \times \underset{\substack{\text{Action} \\ \text{space}}}{A} \rightarrow \mathbb{R}_+ \cup \{+\infty\} =: \overline{\mathbb{R}}$$

- ▶ $a \mapsto \ell(y, a)$ is an outcome contingent utility

- ▶ Learning goal: **expected** risk minimisation

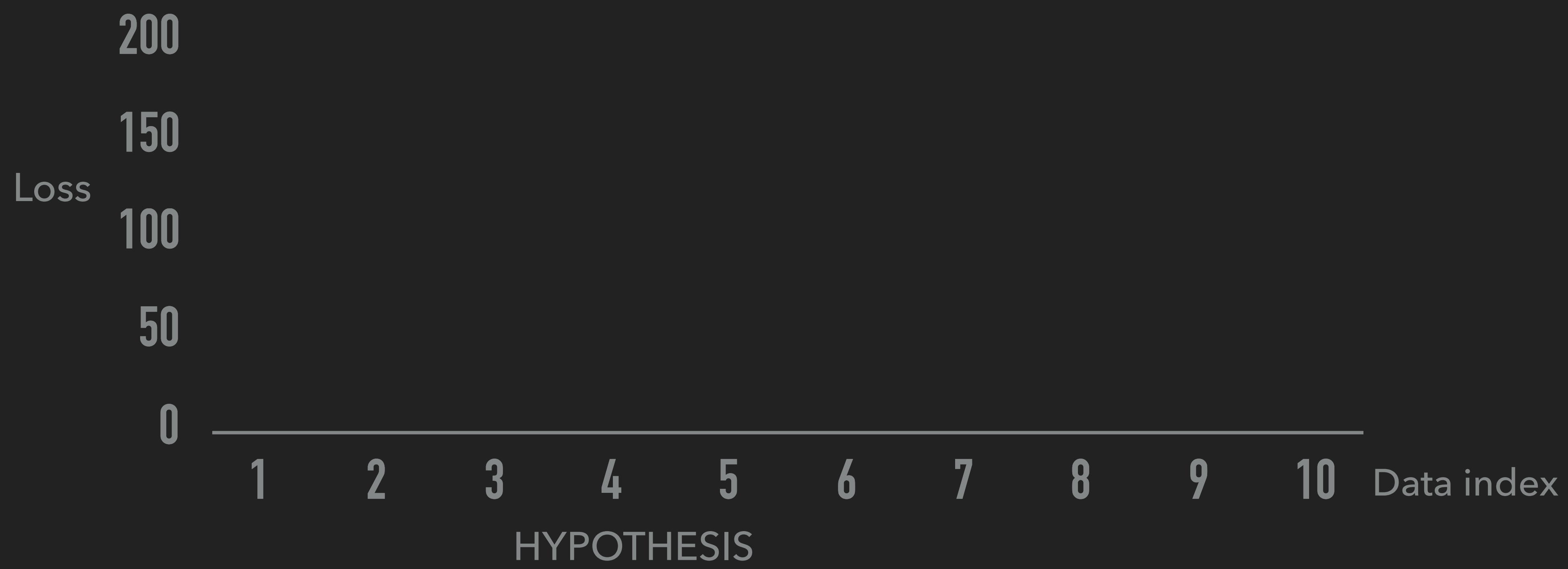
$$\min_{f \in \mathcal{F}} \mathbb{E}_{(X, Y) \sim P} \ell(Y, f(X))$$

- ▶ In practice: **empirical** risk minimisation

$$\begin{aligned} & \min_{f \in \mathcal{F}} \mathbb{E}_{(X, Y) \sim P^m} \ell(Y, f(X)) \\ &= \min_{f \in \mathcal{F}} \frac{1}{m} \sum_{i=1}^m \ell(y_i, f(x_i)) \end{aligned}$$

MINIMISING EMPIRICAL RISK

MINIMISING EMPIRICAL RISK



MINIMISING EMPIRICAL RISK



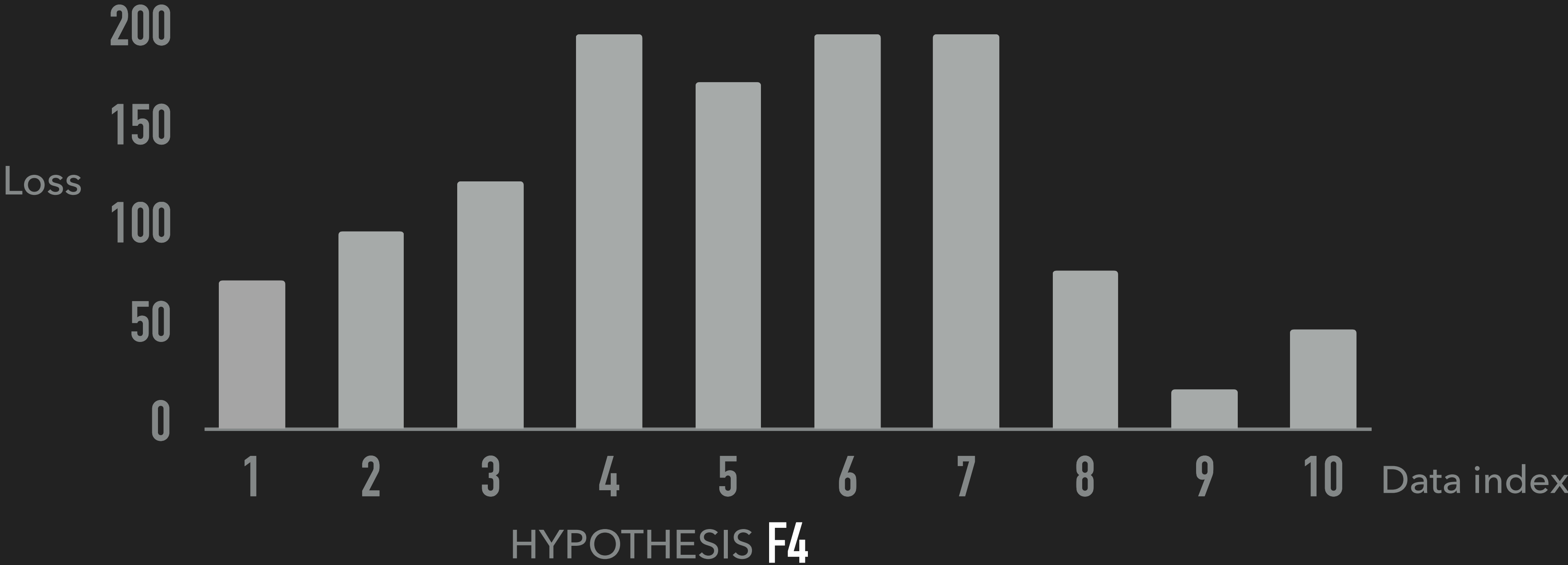
MINIMISING EMPIRICAL RISK



MINIMISING EMPIRICAL RISK



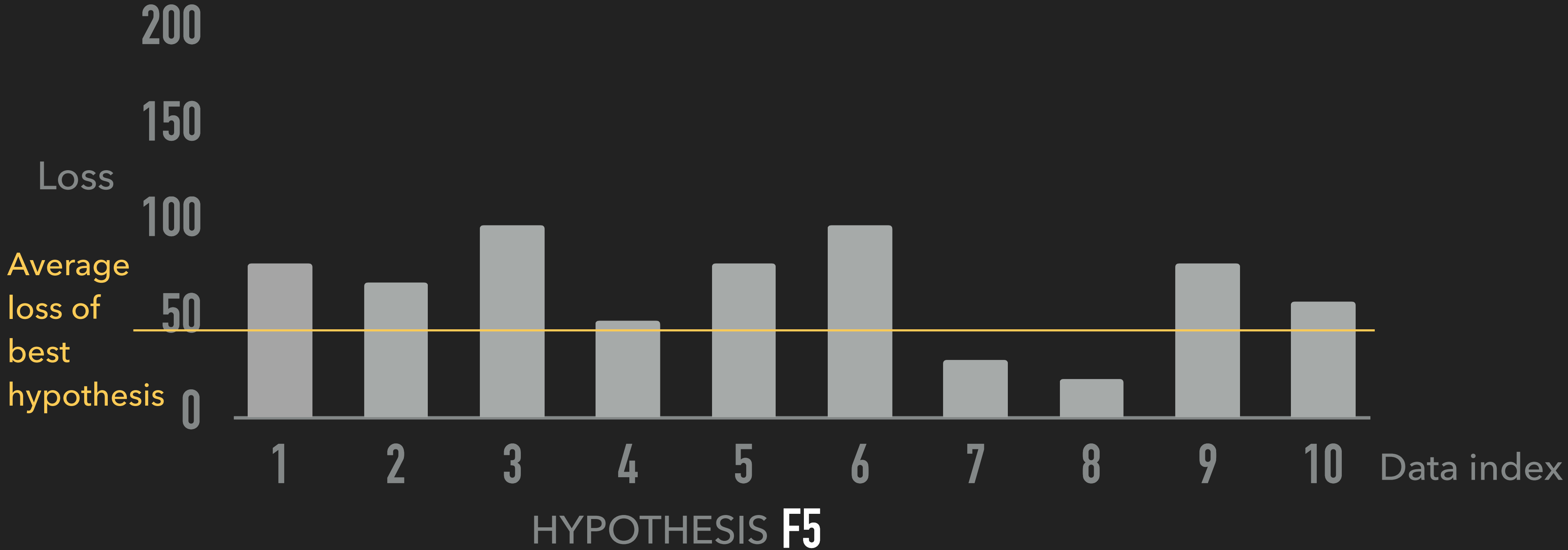
MINIMISING EMPIRICAL RISK



MINIMISING EMPIRICAL RISK



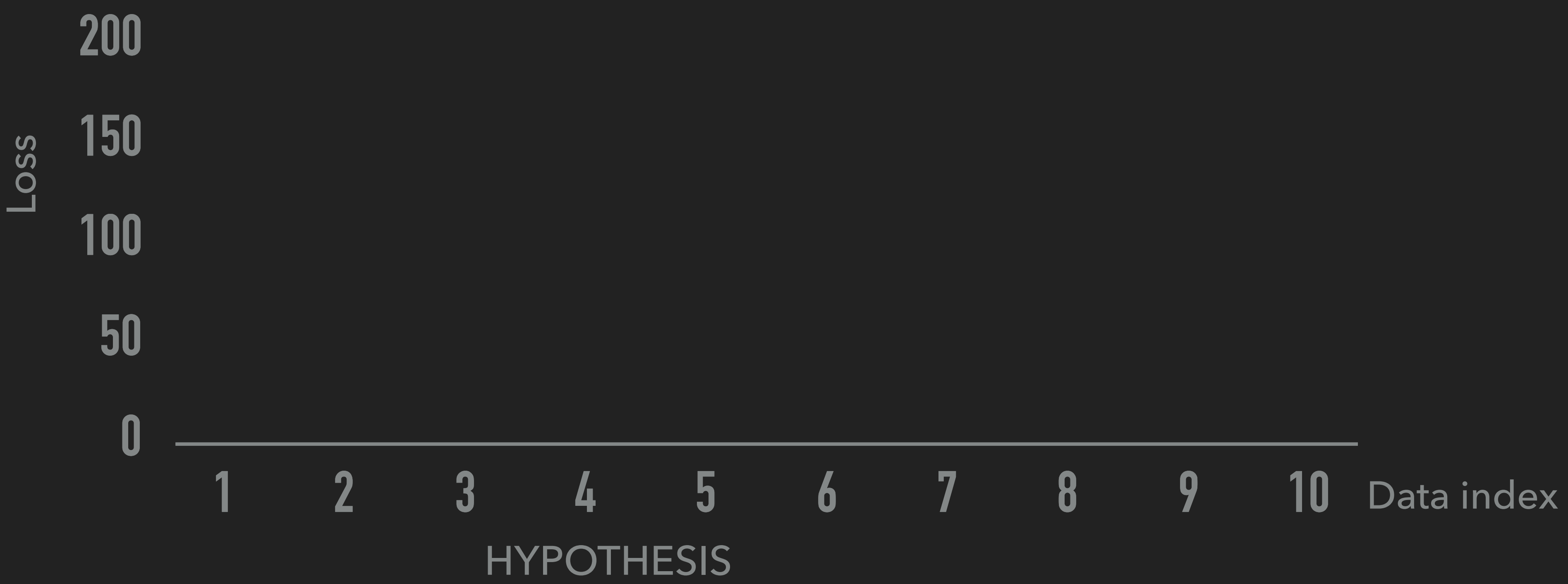
MINIMISING EMPIRICAL RISK



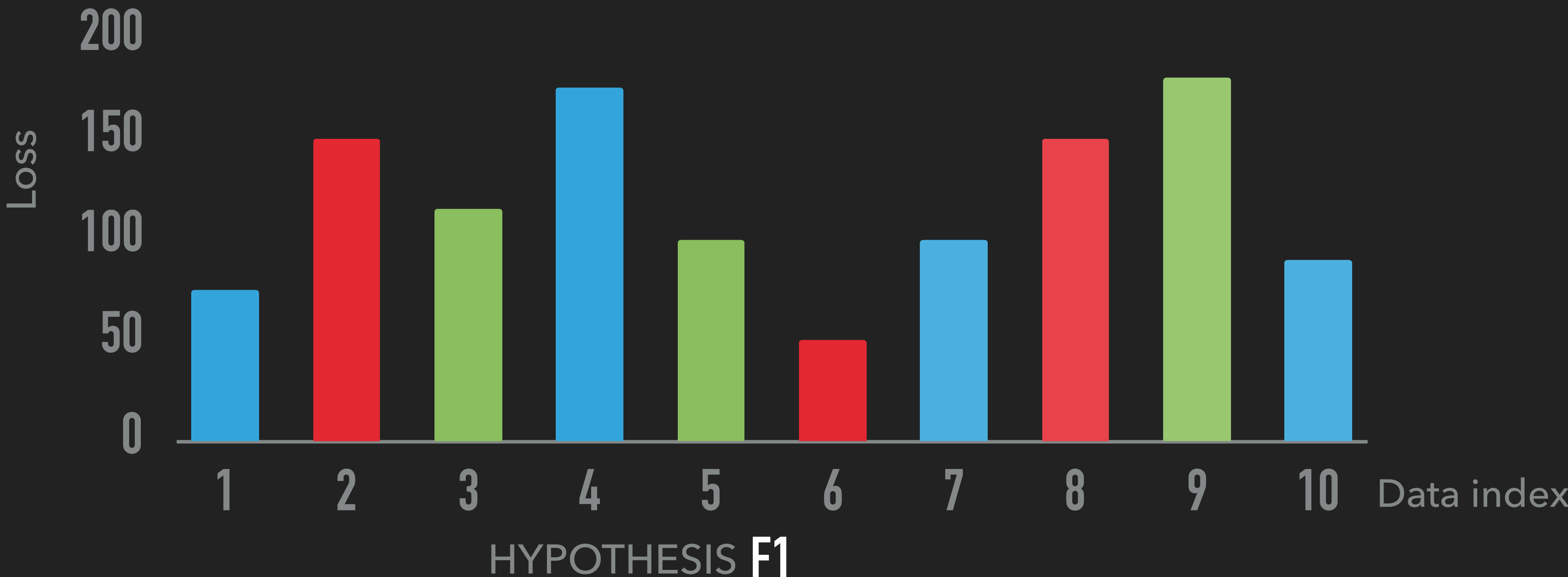
HYPOTHESIS **F5**

MINIMISING EMPIRICAL RISK WITH SENSITIVE ATTRIBUTES VISIBLE

MINIMISING EMPIRICAL RISK WITH SENSITIVE ATTRIBUTES VISIBLE

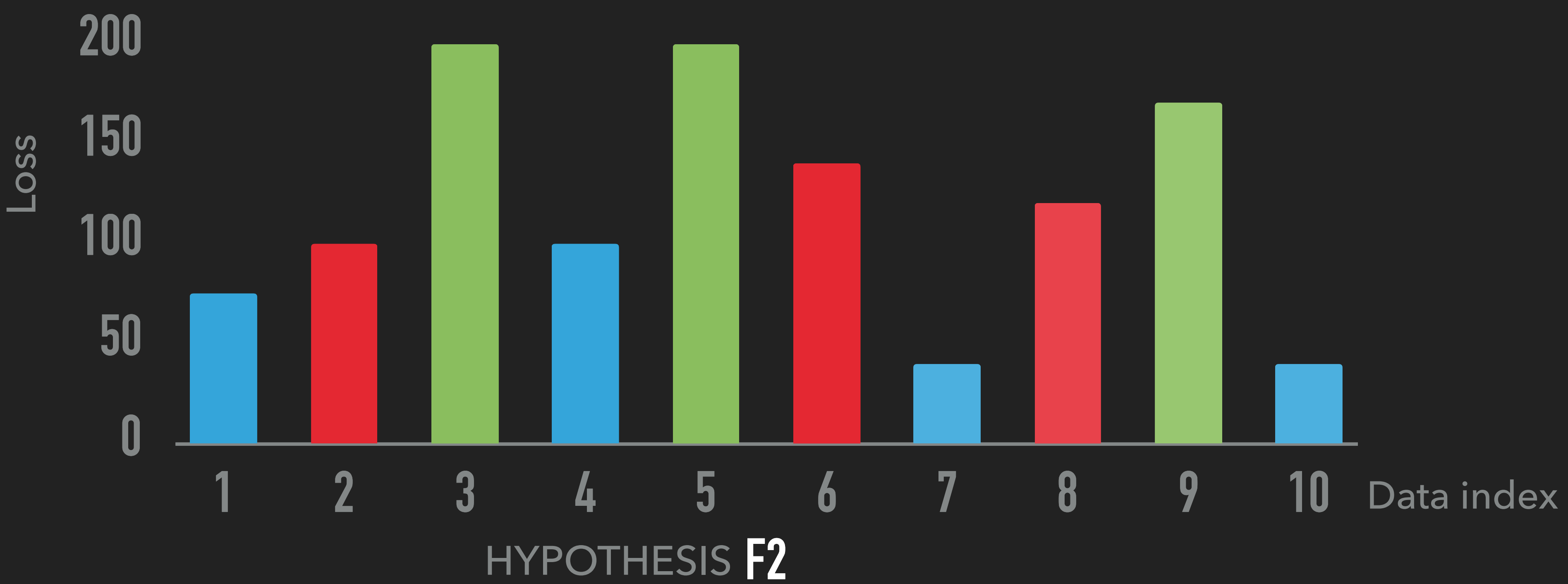


MINIMISING EMPIRICAL RISK WITH SENSITIVE ATTRIBUTES VISIBLE

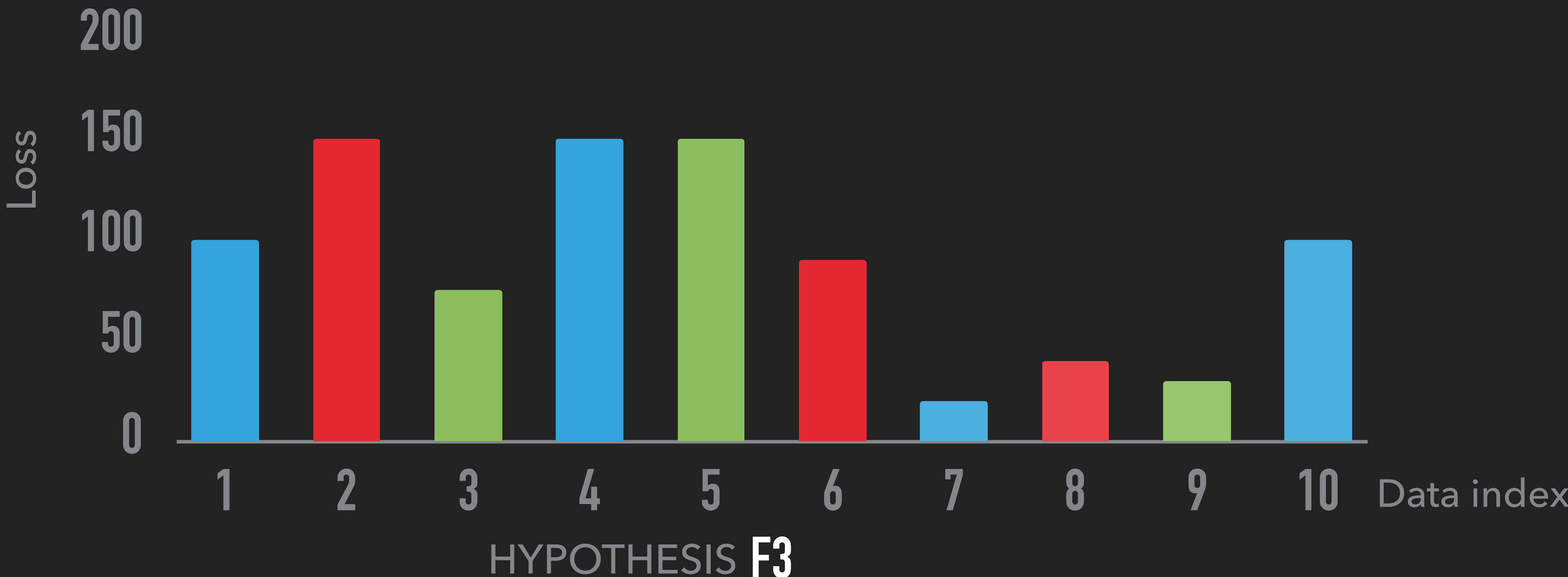


HYPOTHESIS **F1**

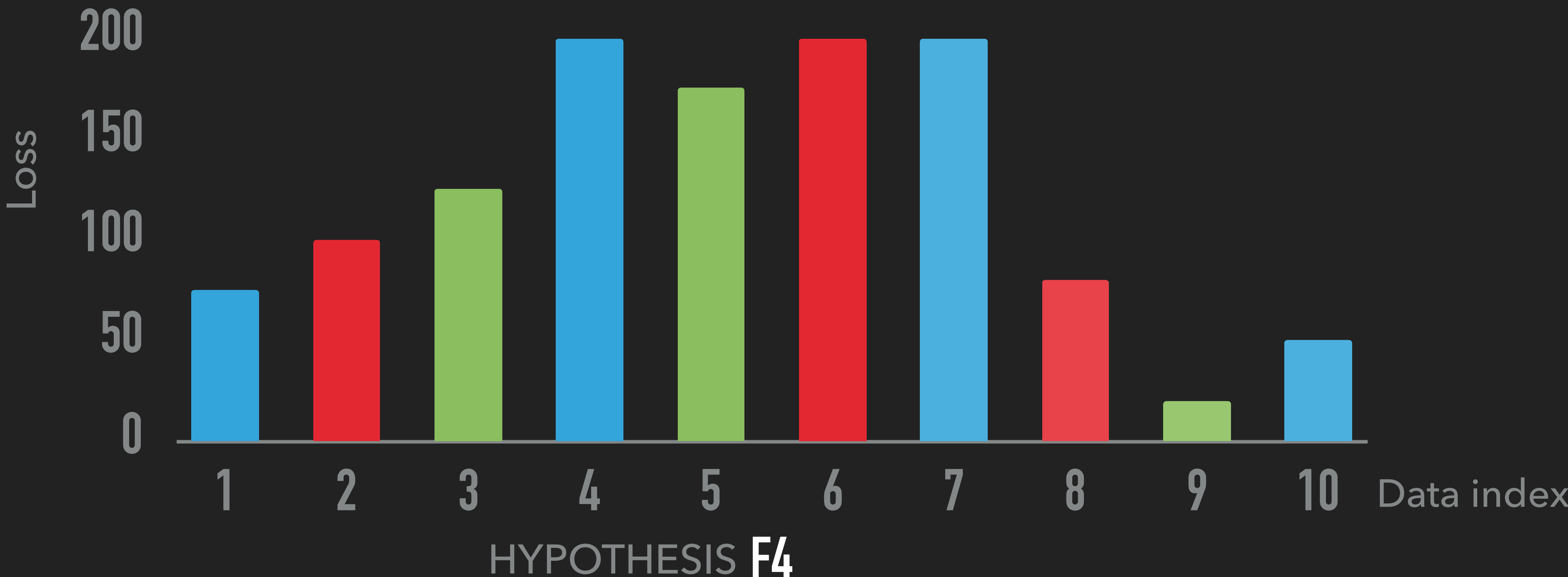
MINIMISING EMPIRICAL RISK WITH SENSITIVE ATTRIBUTES VISIBLE



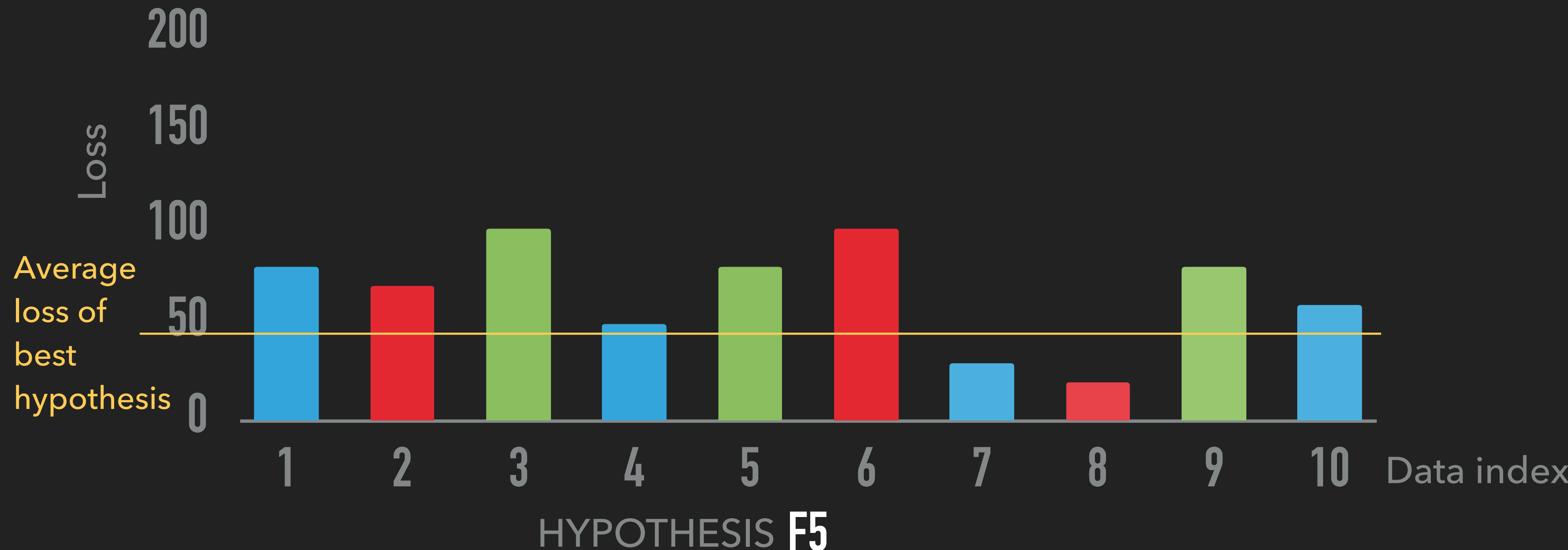
MINIMISING EMPIRICAL RISK WITH SENSITIVE ATTRIBUTES VISIBLE



MINIMISING EMPIRICAL RISK WITH SENSITIVE ATTRIBUTES VISIBLE



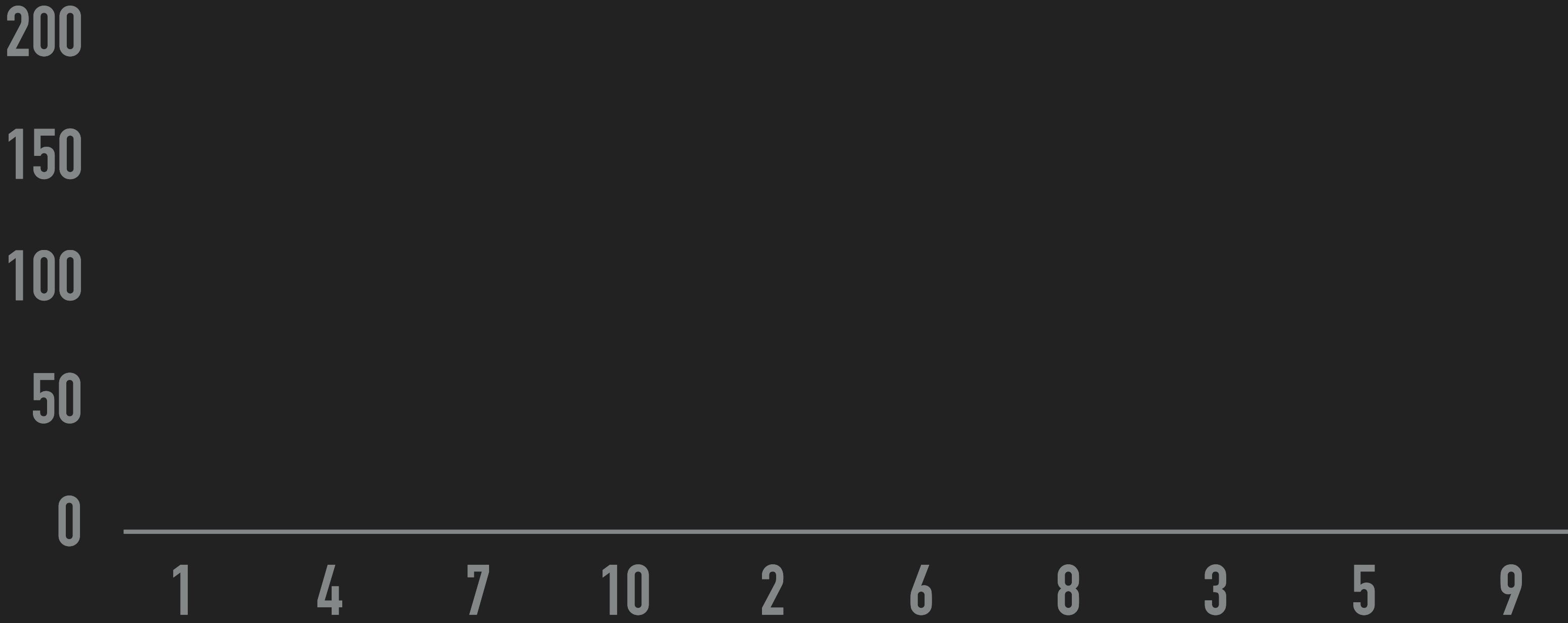
MINIMISING EMPIRICAL RISK WITH SENSITIVE ATTRIBUTES VISIBLE



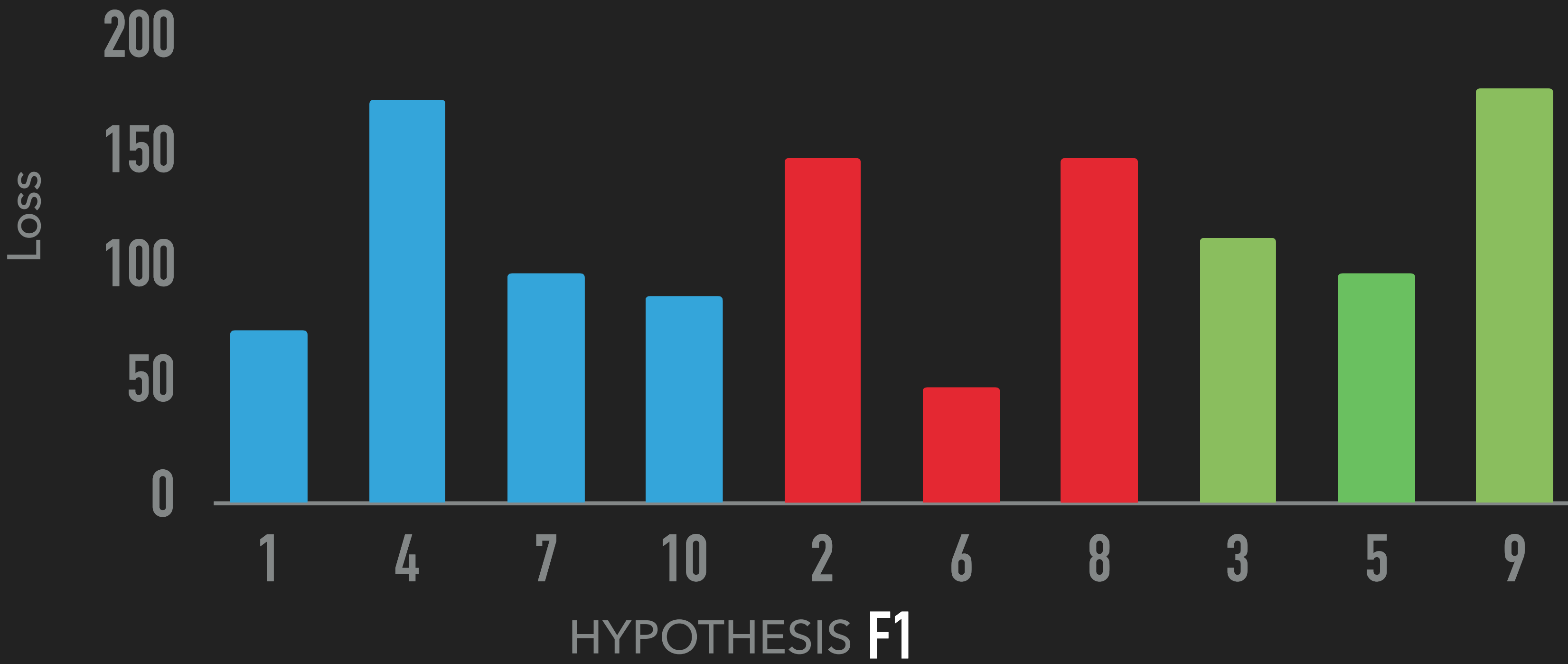
HYPOTHESIS F5

MINIMISING EMPIRICAL RISK REINDEXING PER SENSITIVE ATTRIBUTE

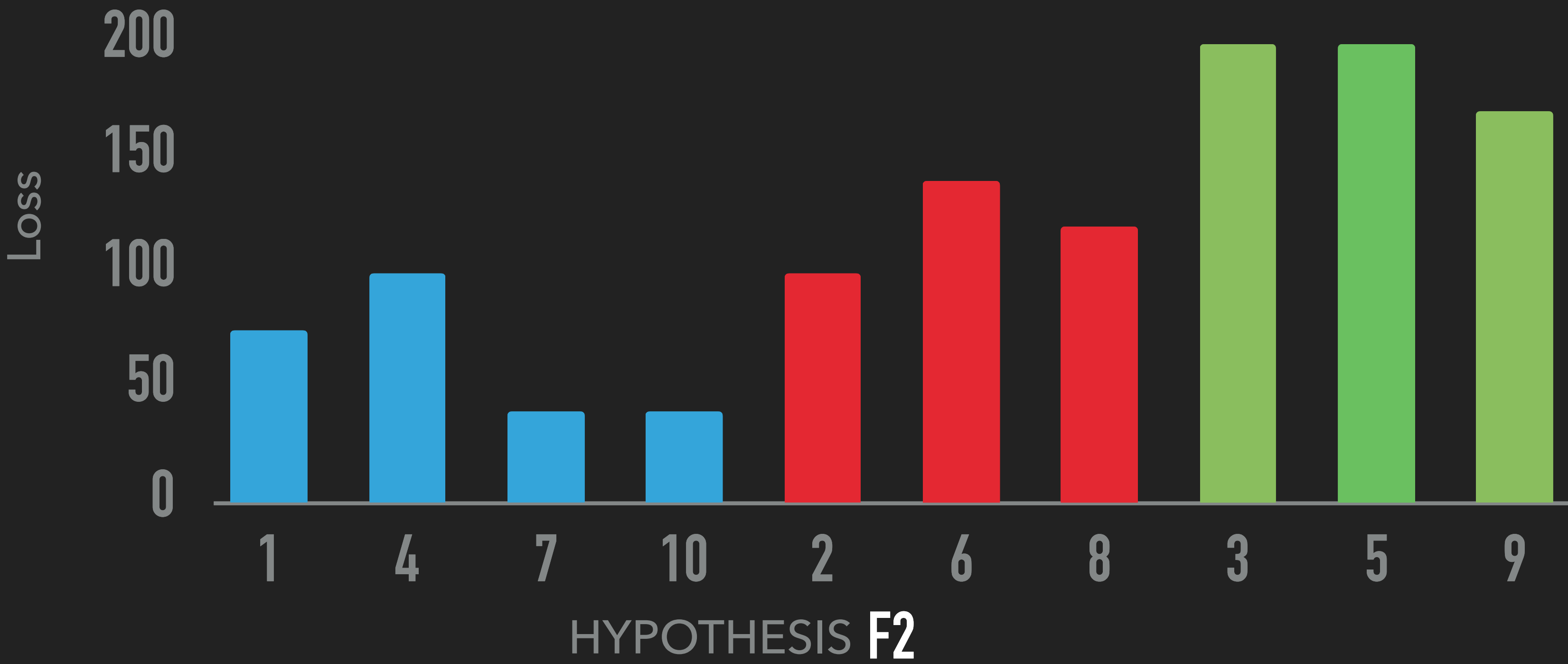
MINIMISING EMPIRICAL RISK REINDEXING PER SENSITIVE ATTRIBUTE



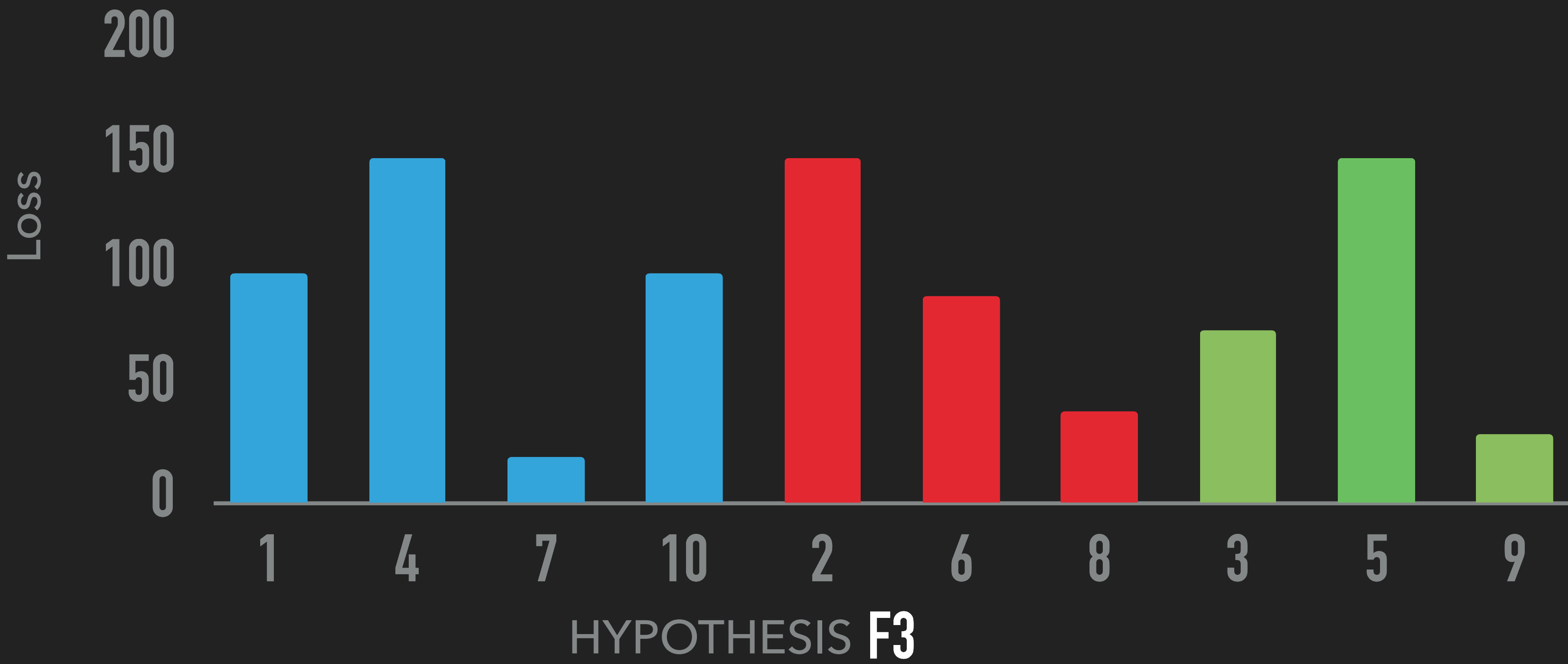
MINIMISING EMPIRICAL RISK REINDEXING PER SENSITIVE ATTRIBUTE



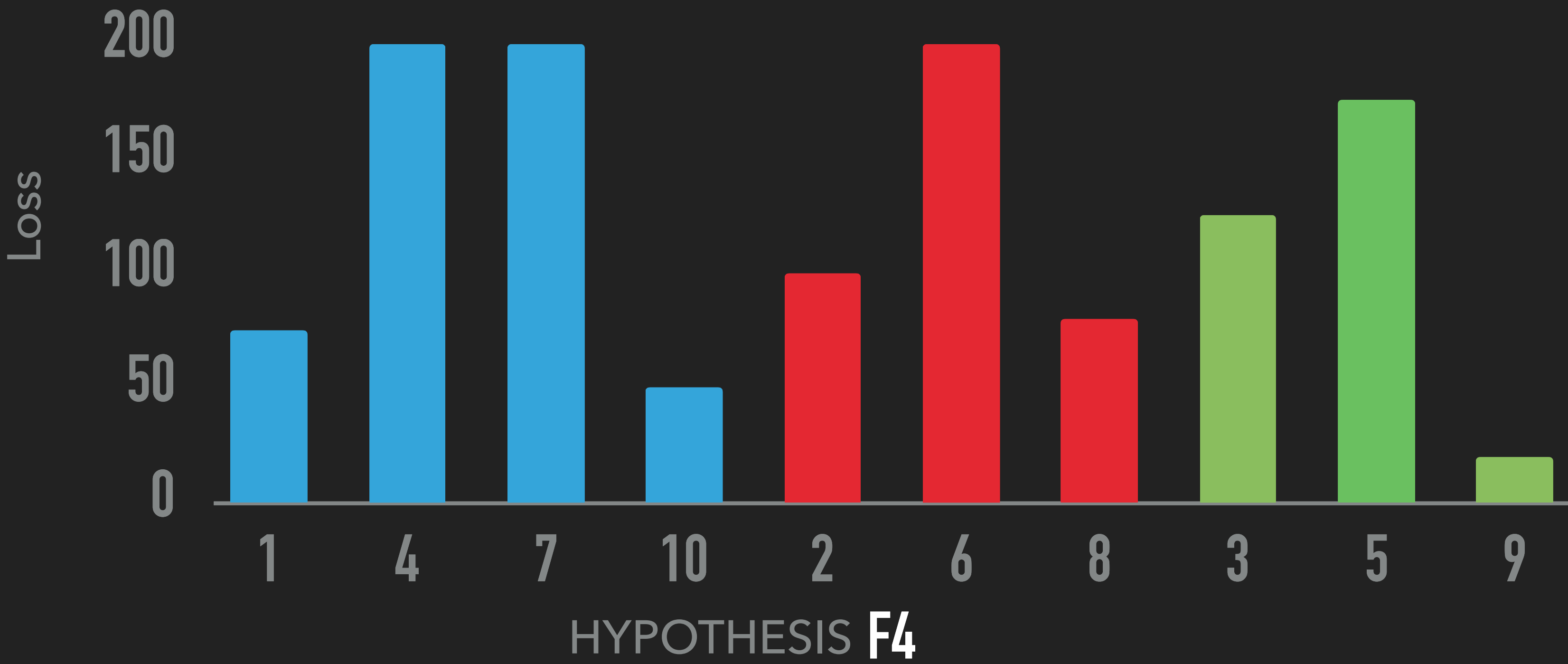
MINIMISING EMPIRICAL RISK REINDEXING PER SENSITIVE ATTRIBUTE



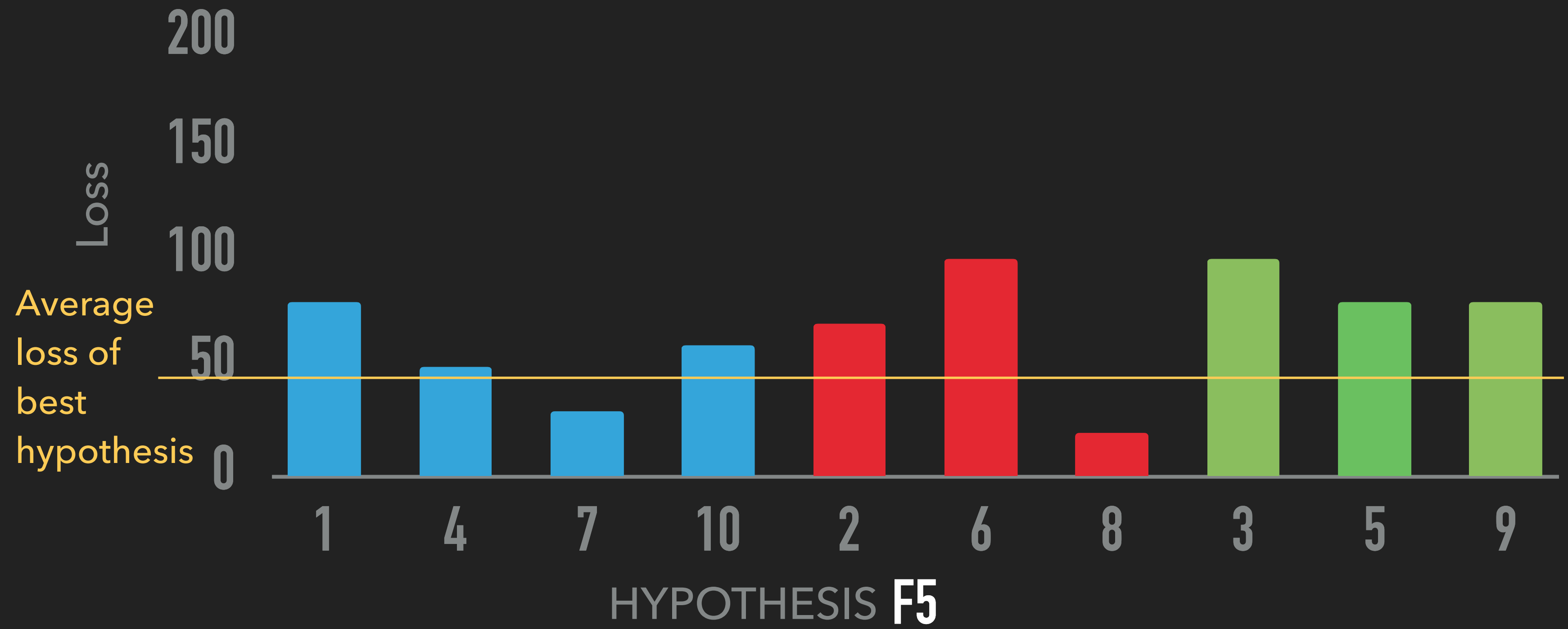
MINIMISING EMPIRICAL RISK REINDEXING PER SENSITIVE ATTRIBUTE



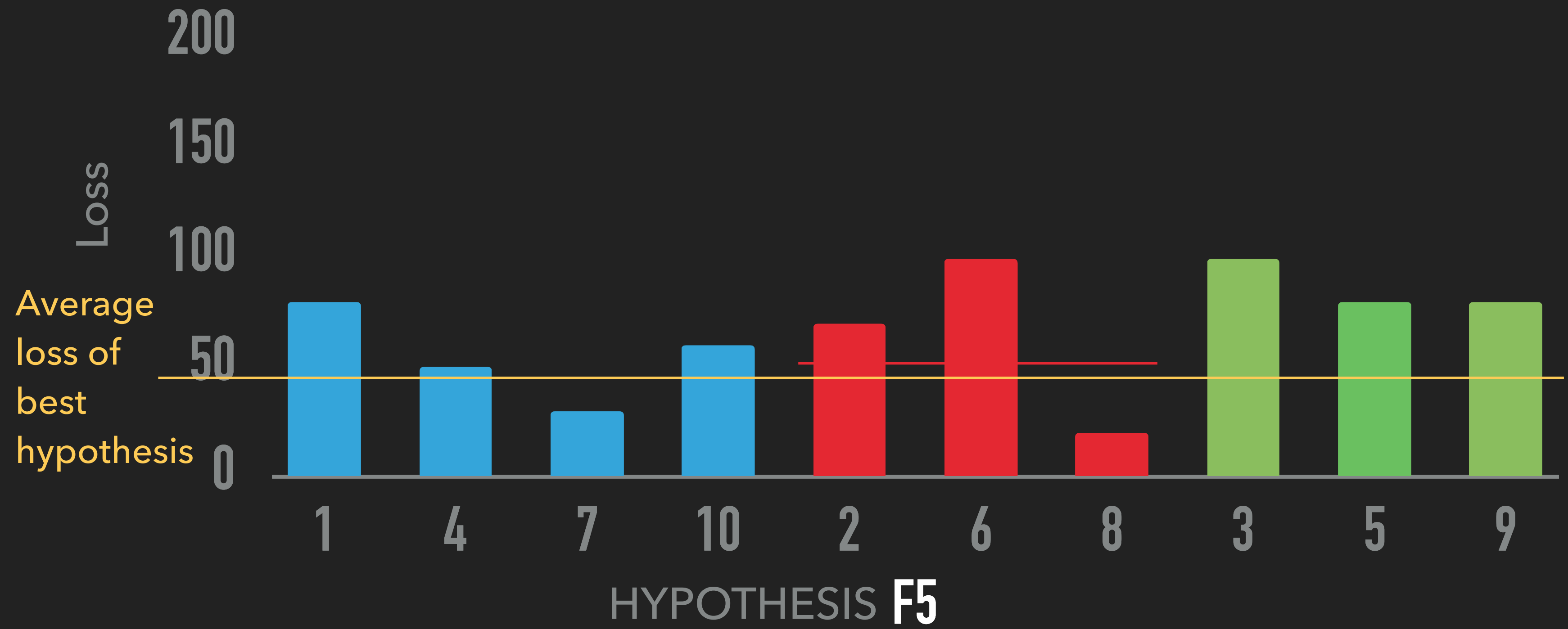
MINIMISING EMPIRICAL RISK REINDEXING PER SENSITIVE ATTRIBUTE



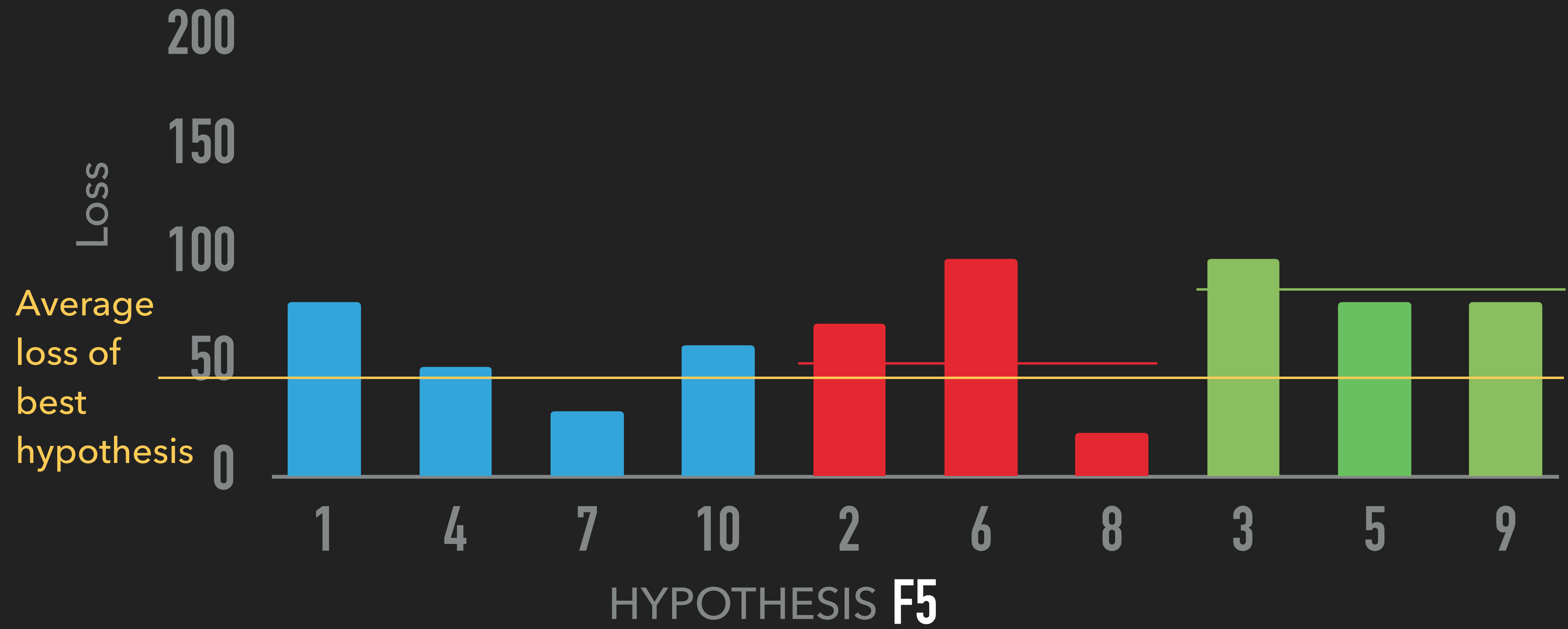
MINIMISING EMPIRICAL RISK REINDEXING PER SENSITIVE ATTRIBUTE



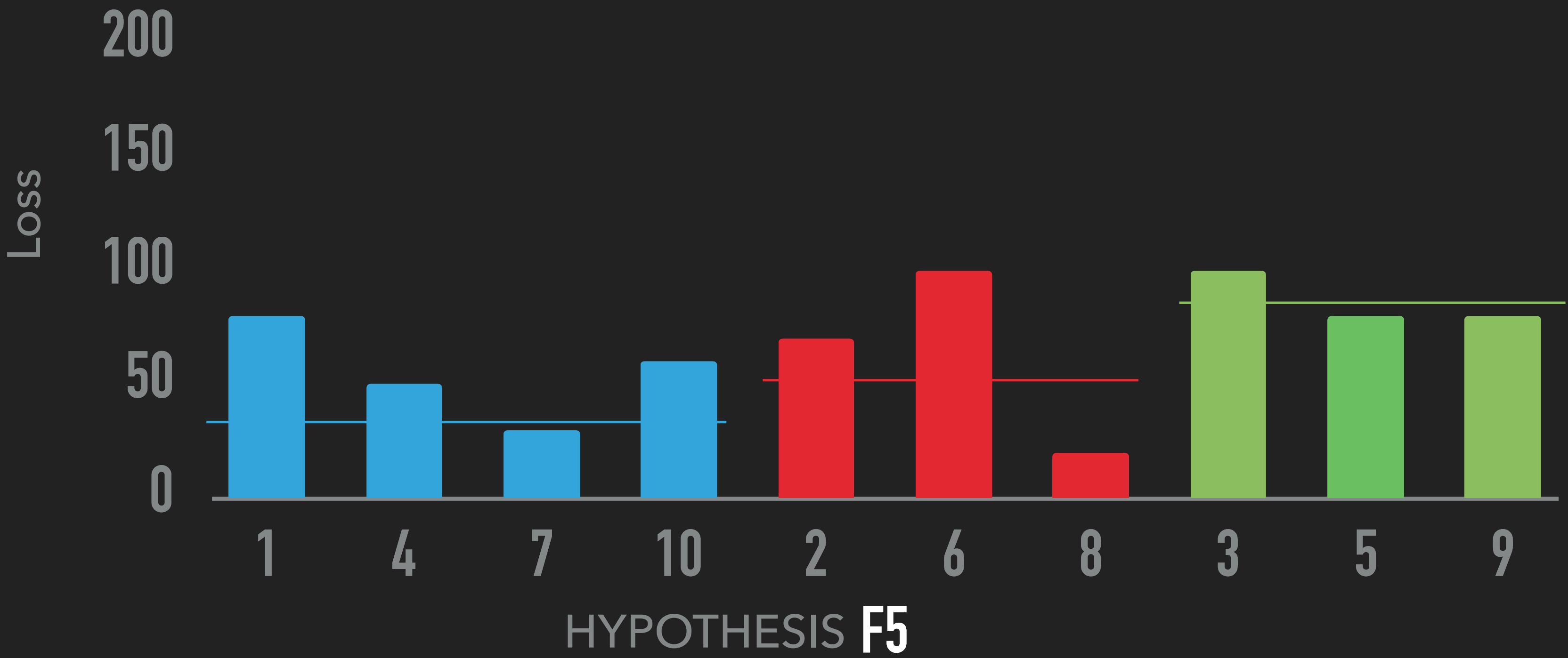
MINIMISING EMPIRICAL RISK REINDEXING PER SENSITIVE ATTRIBUTE



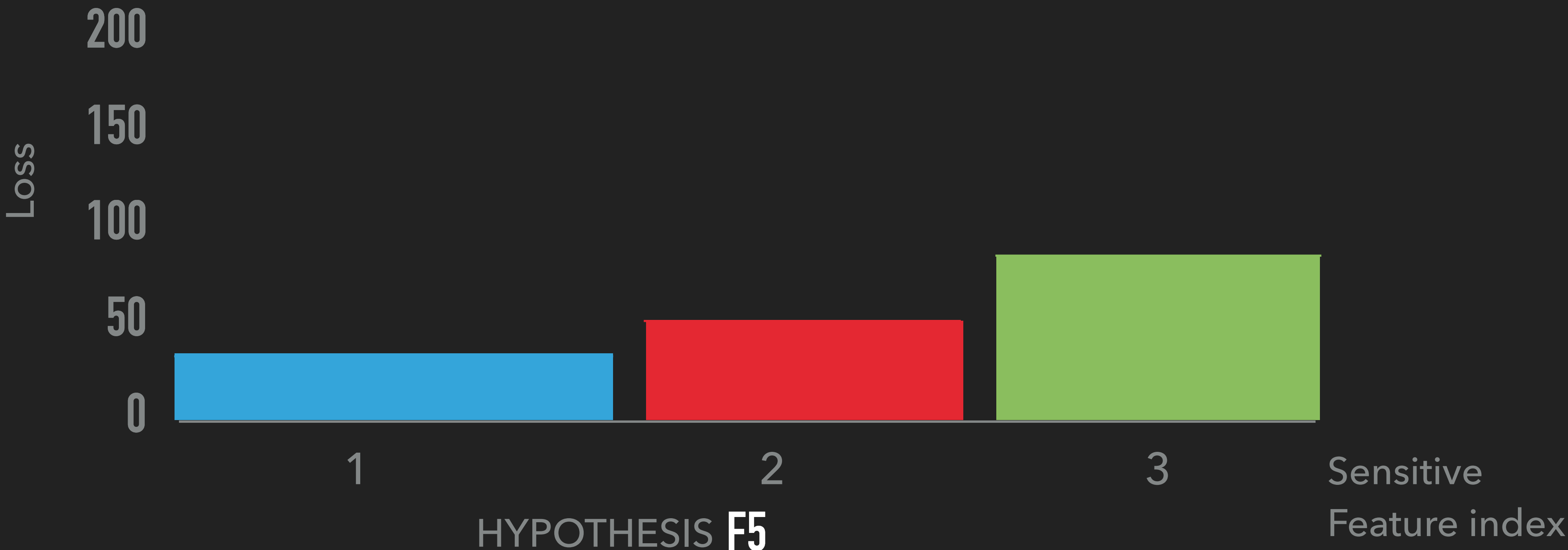
MINIMISING EMPIRICAL RISK REINDEXING PER SENSITIVE ATTRIBUTE



MINIMISING EMPIRICAL RISK REINDEXING PER SENSITIVE ATTRIBUTE

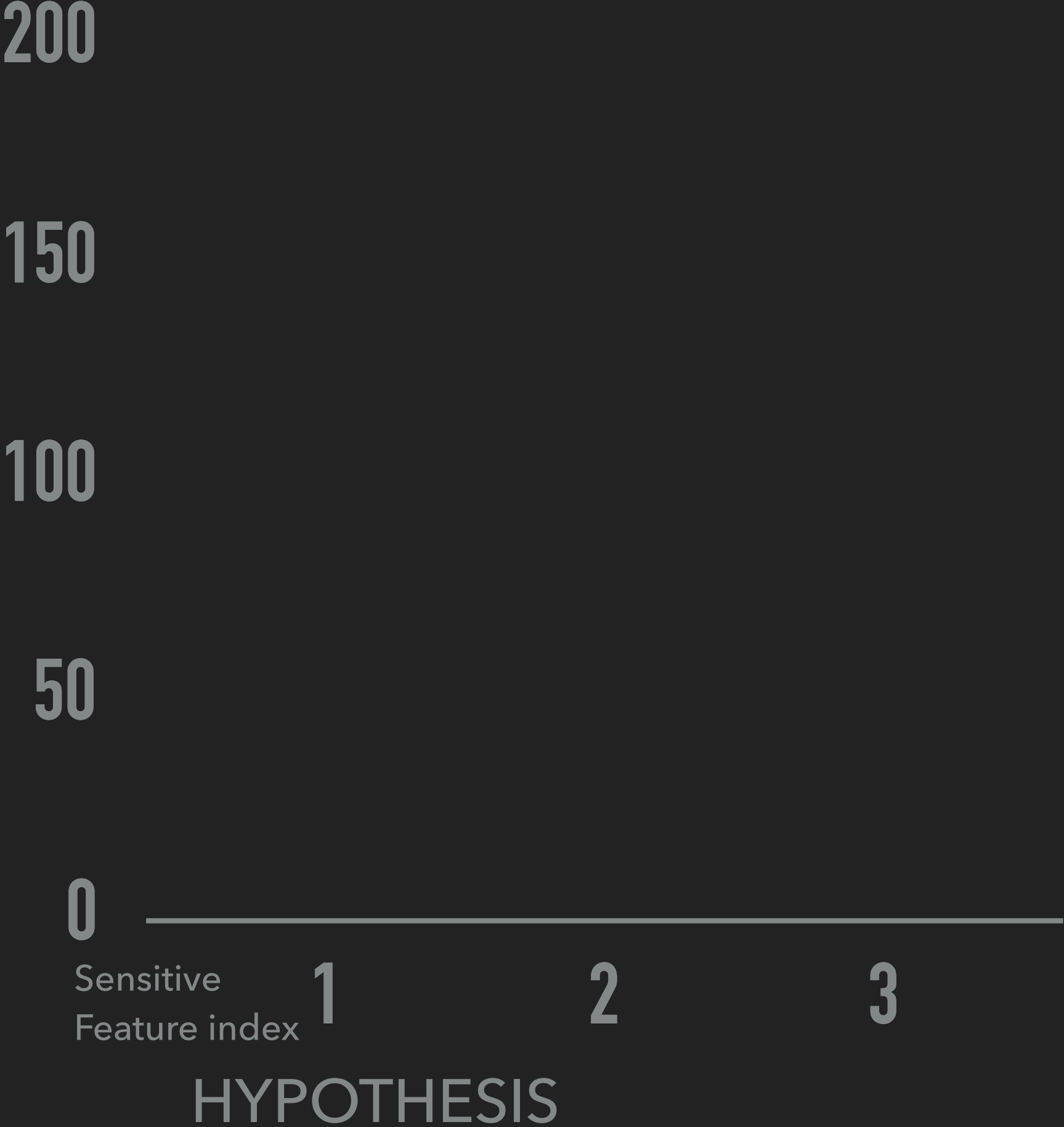


MINIMISING EMPIRICAL RISK REINDEXING PER SENSITIVE ATTRIBUTE

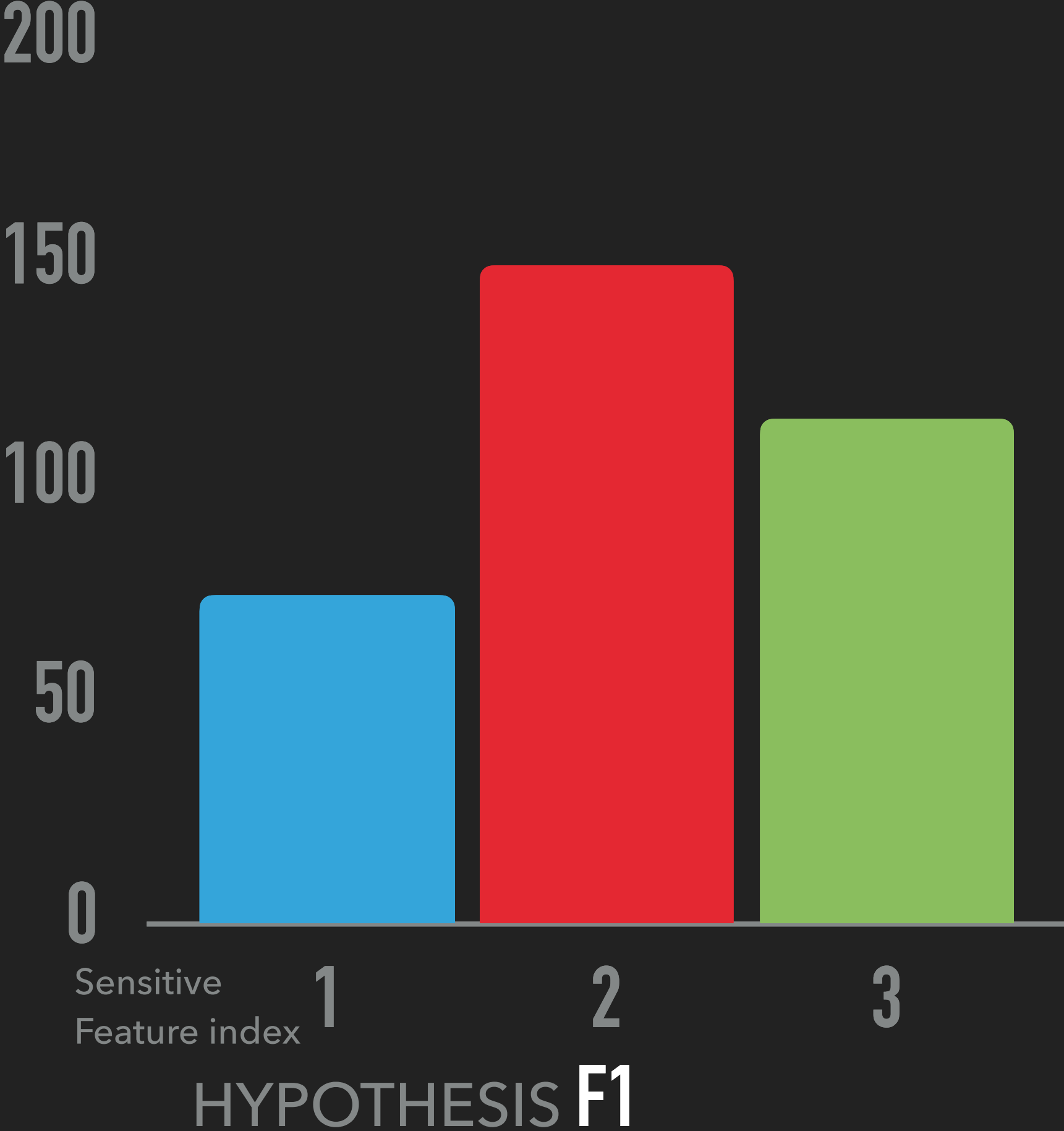


MINIMISING AGGREGATED EMPIRICAL RISK

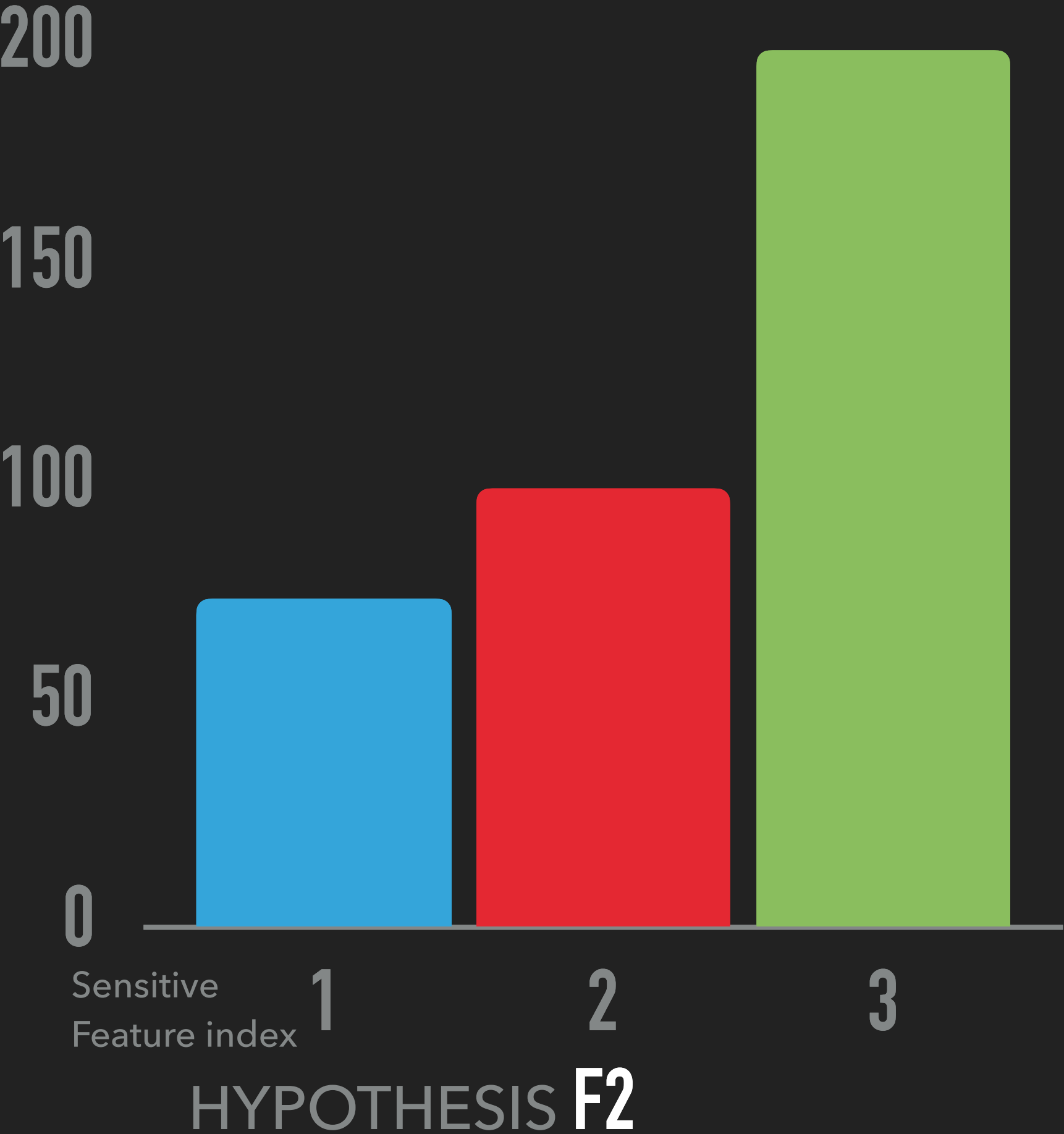
MINIMISING AGGREGATED EMPIRICAL RISK



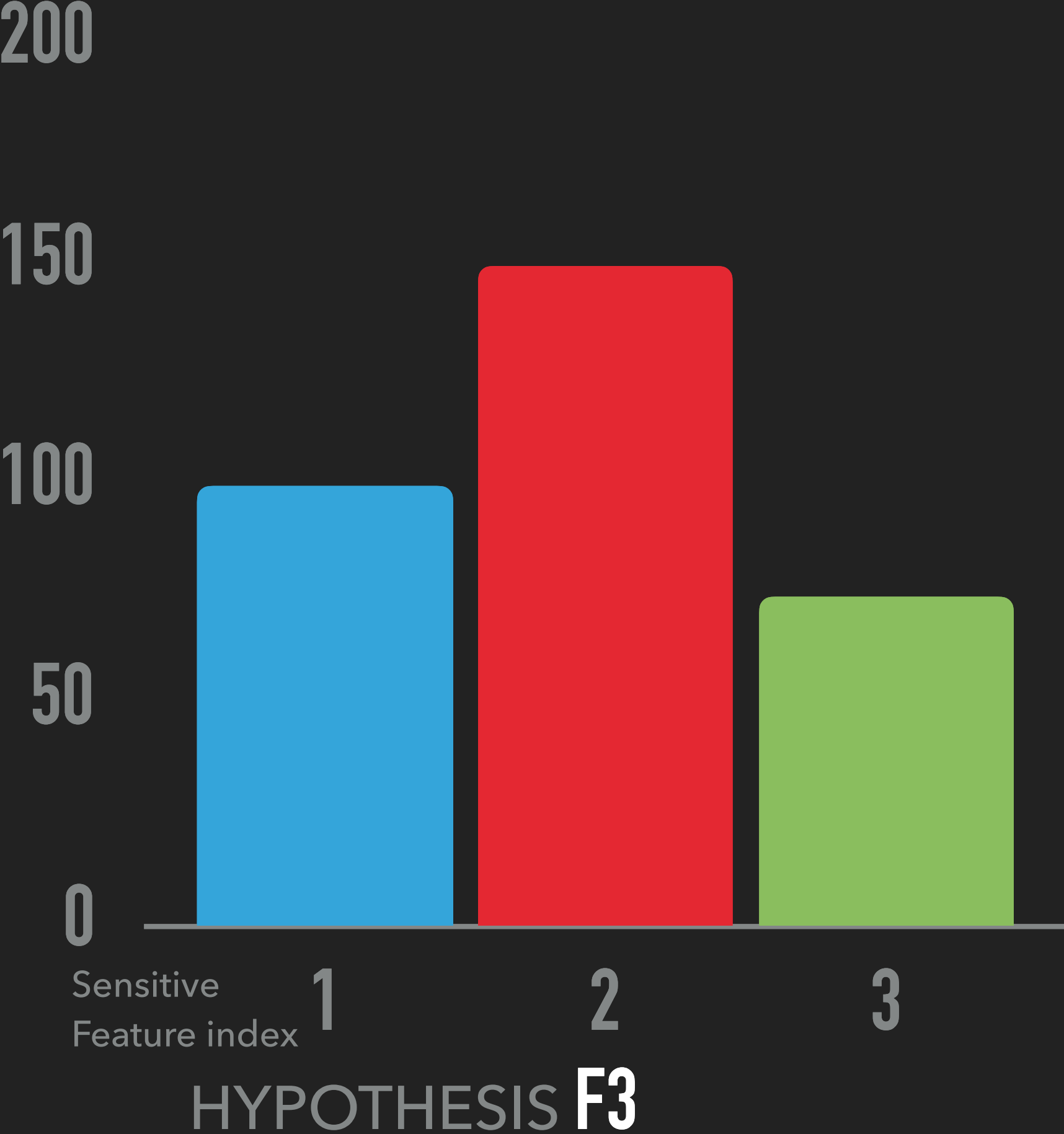
MINIMISING AGGREGATED EMPIRICAL RISK



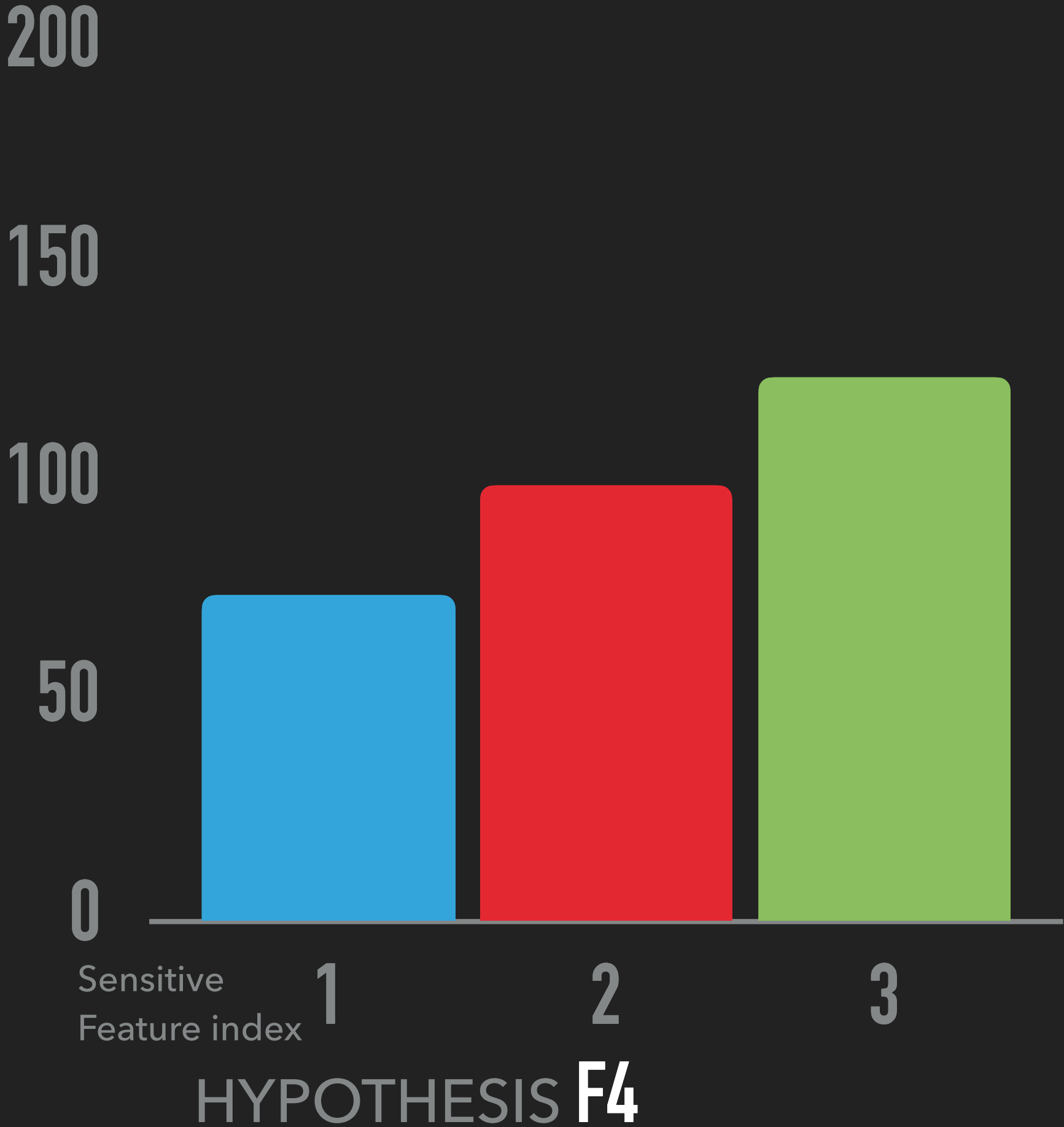
MINIMISING AGGREGATED EMPIRICAL RISK



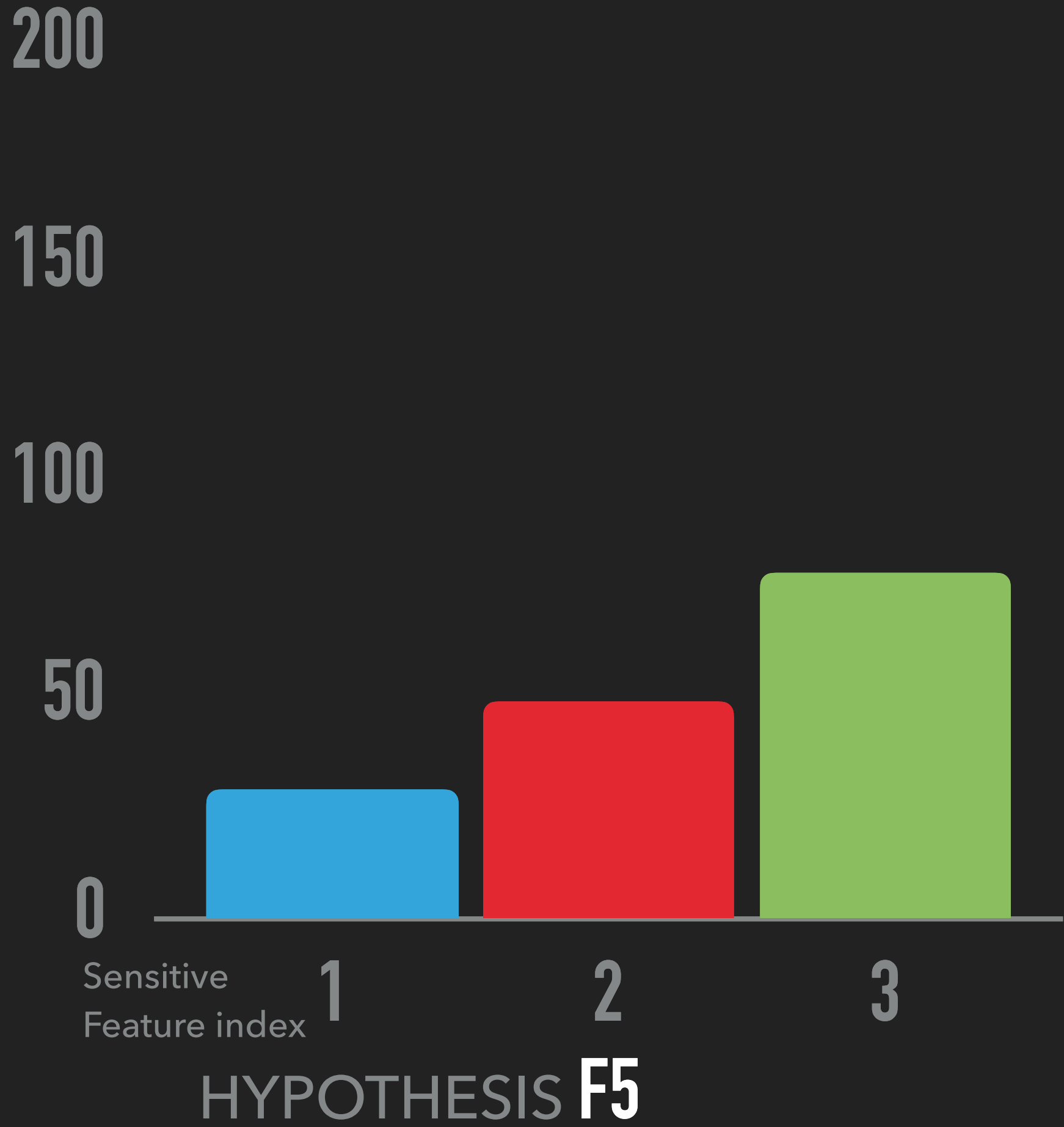
MINIMISING AGGREGATED EMPIRICAL RISK



MINIMISING AGGREGATED EMPIRICAL RISK

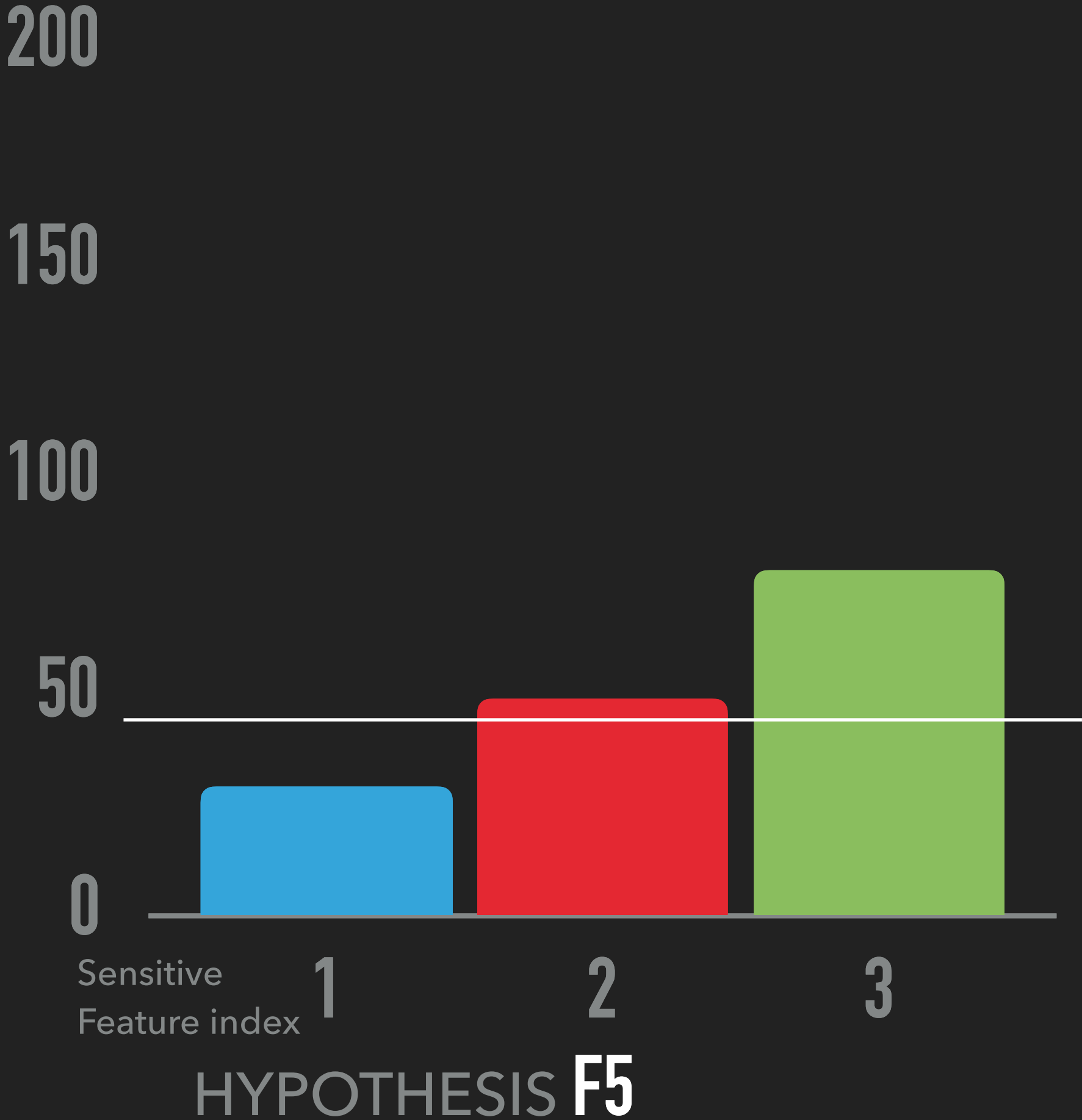


MINIMISING AGGREGATED EMPIRICAL RISK



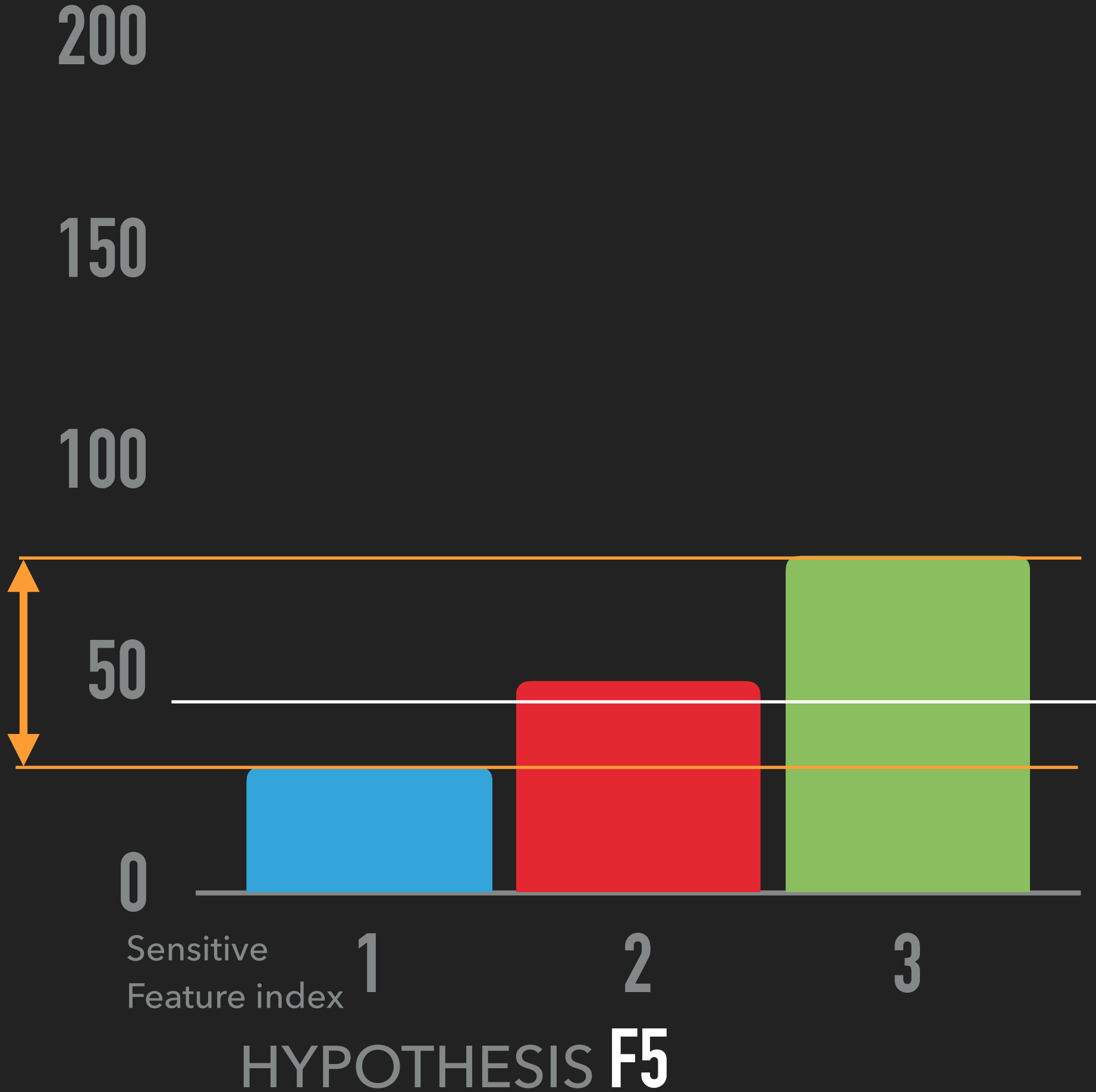
MINIMISING AGGREGATED EMPIRICAL RISK

▶ Standard problem: minimise average risk



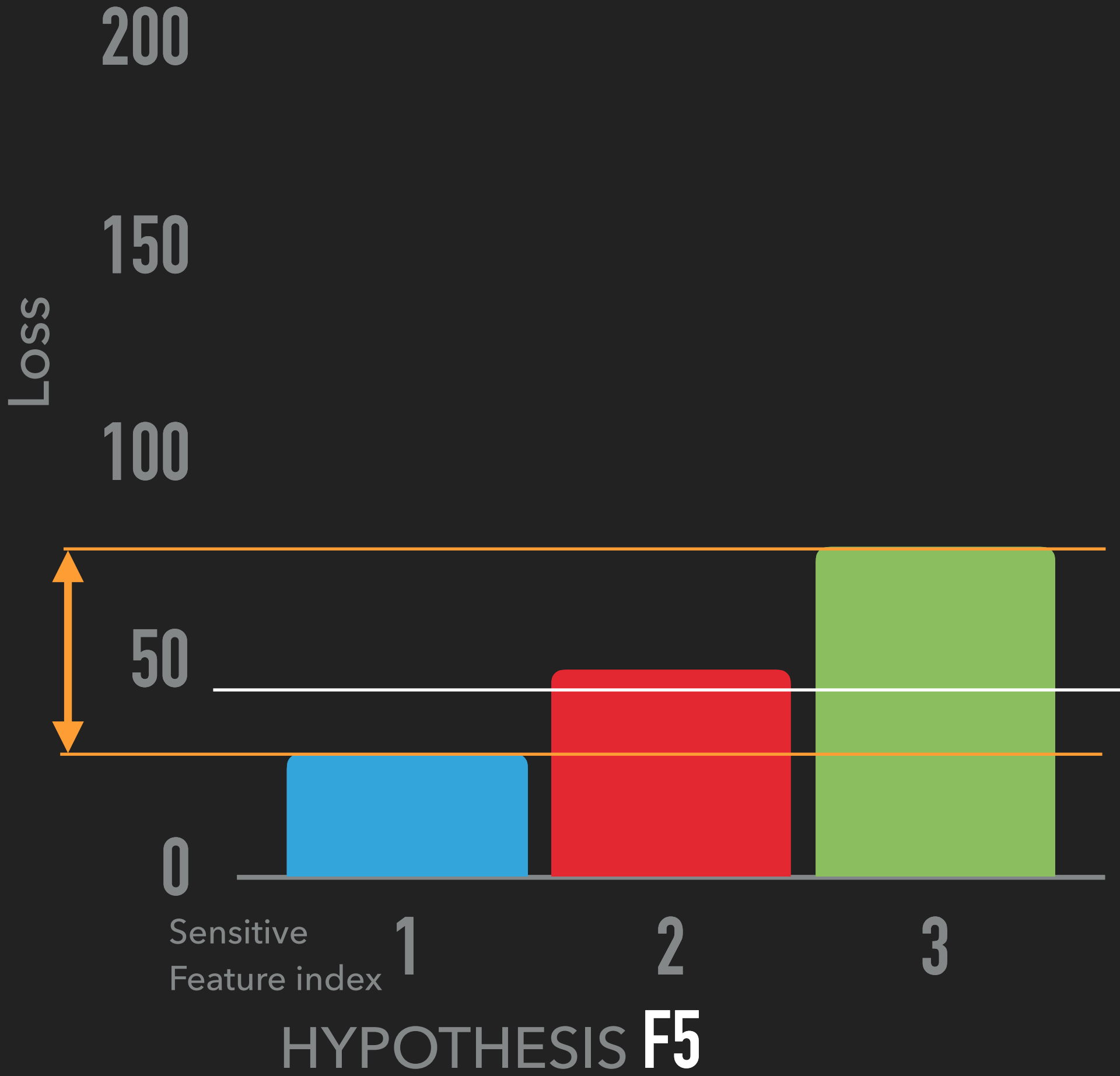
MINIMISING AGGREGATED EMPIRICAL RISK

- ▶ Standard problem: minimise average risk
- ▶ Equity problem: also take account of **variation**



MINIMISING AGGREGATED EMPIRICAL RISK

- ▶ **Standard problem:** minimise average risk
- ▶ **Equity problem:** also take account of **variation**
- ▶ **Fairness problem:** mixture of both



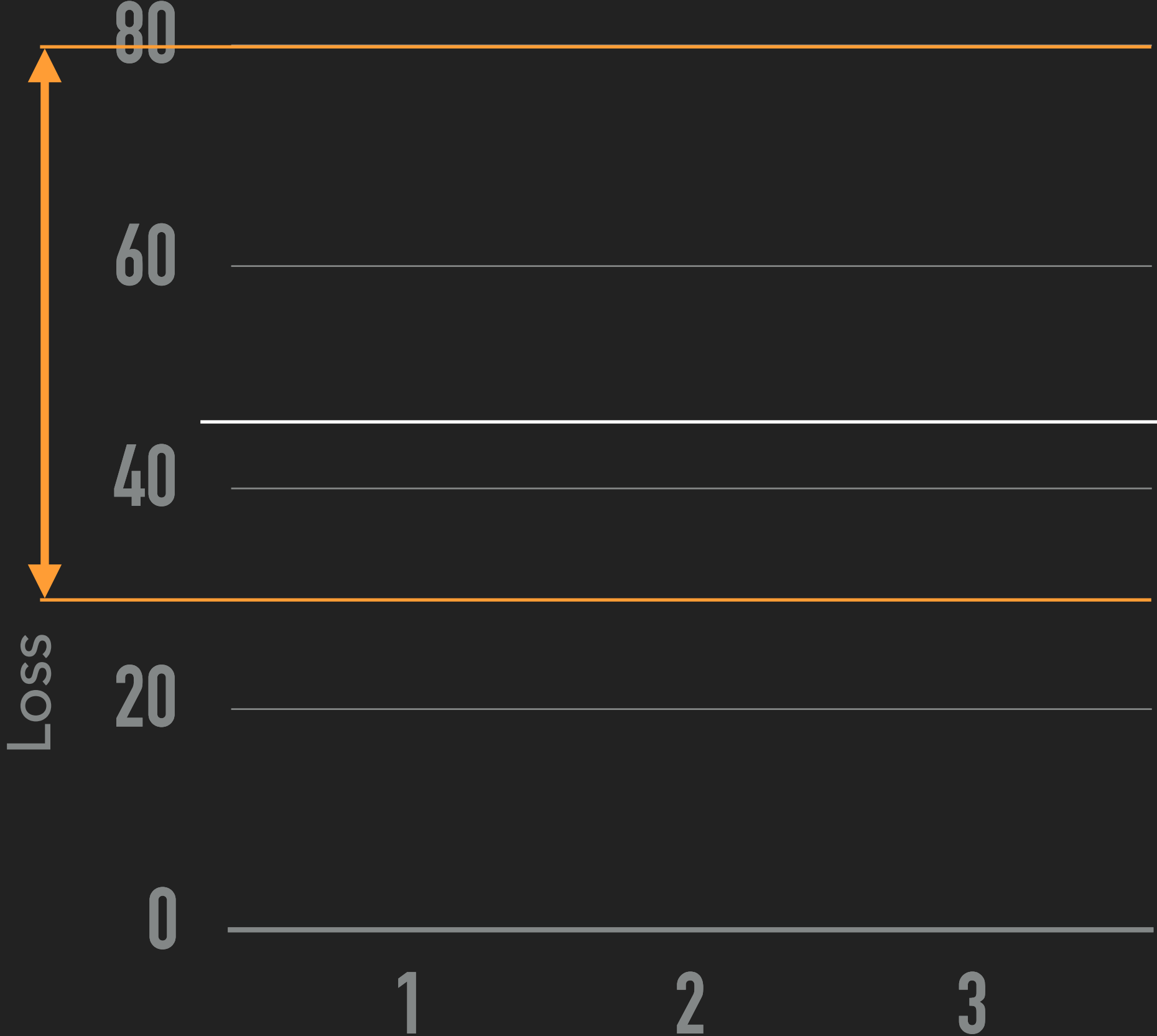
MINIMISING AGGREGATED EMPIRICAL RISK AND DEVIATION

MINIMISING AGGREGATED EMPIRICAL RISK AND DEVIATION

- ▶ Trade off low **deviation** against higher **average**

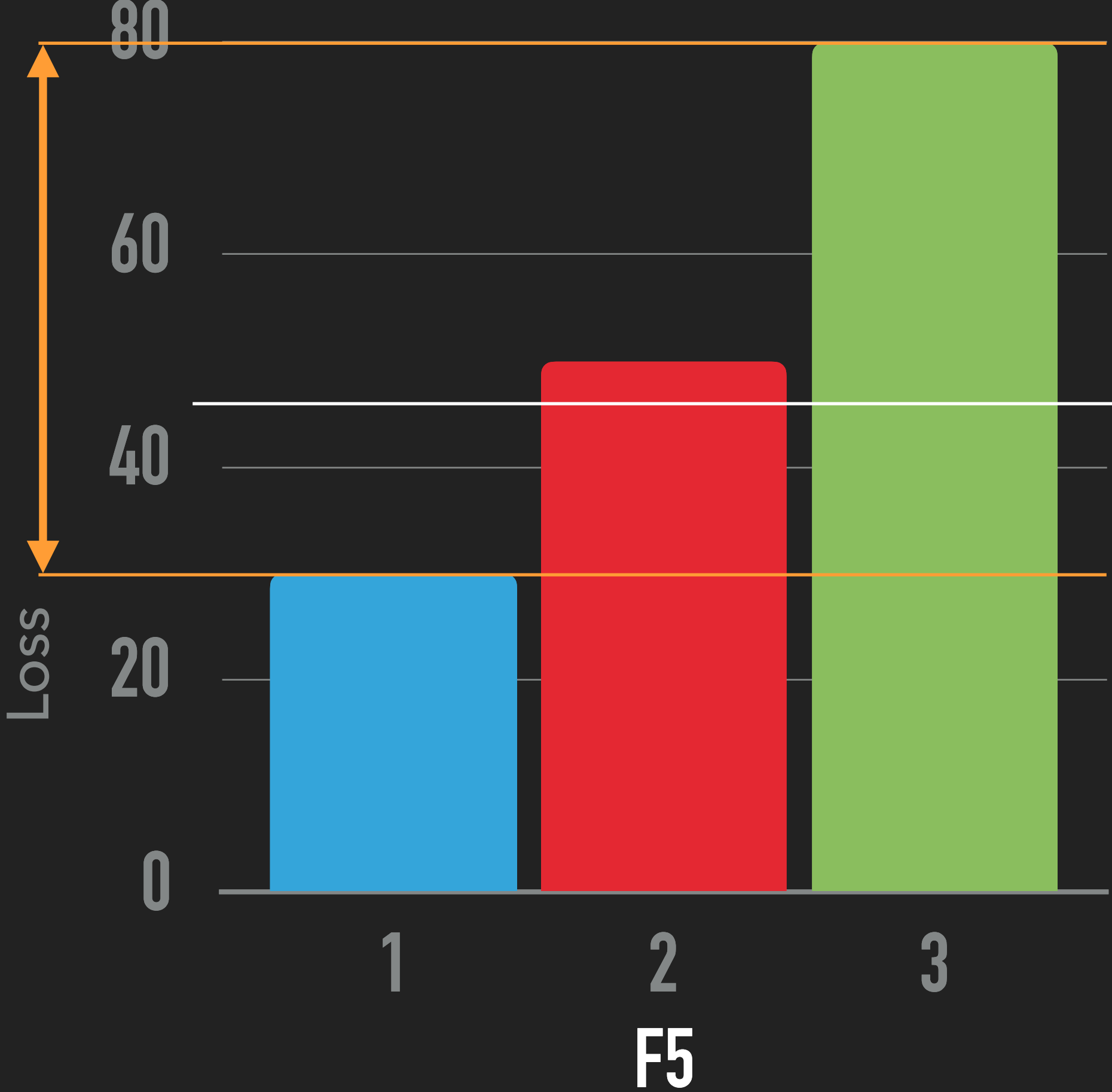
MINIMISING AGGREGATED EMPIRICAL RISK AND DEVIATION

▶ Trade off low **deviation** against higher **average**



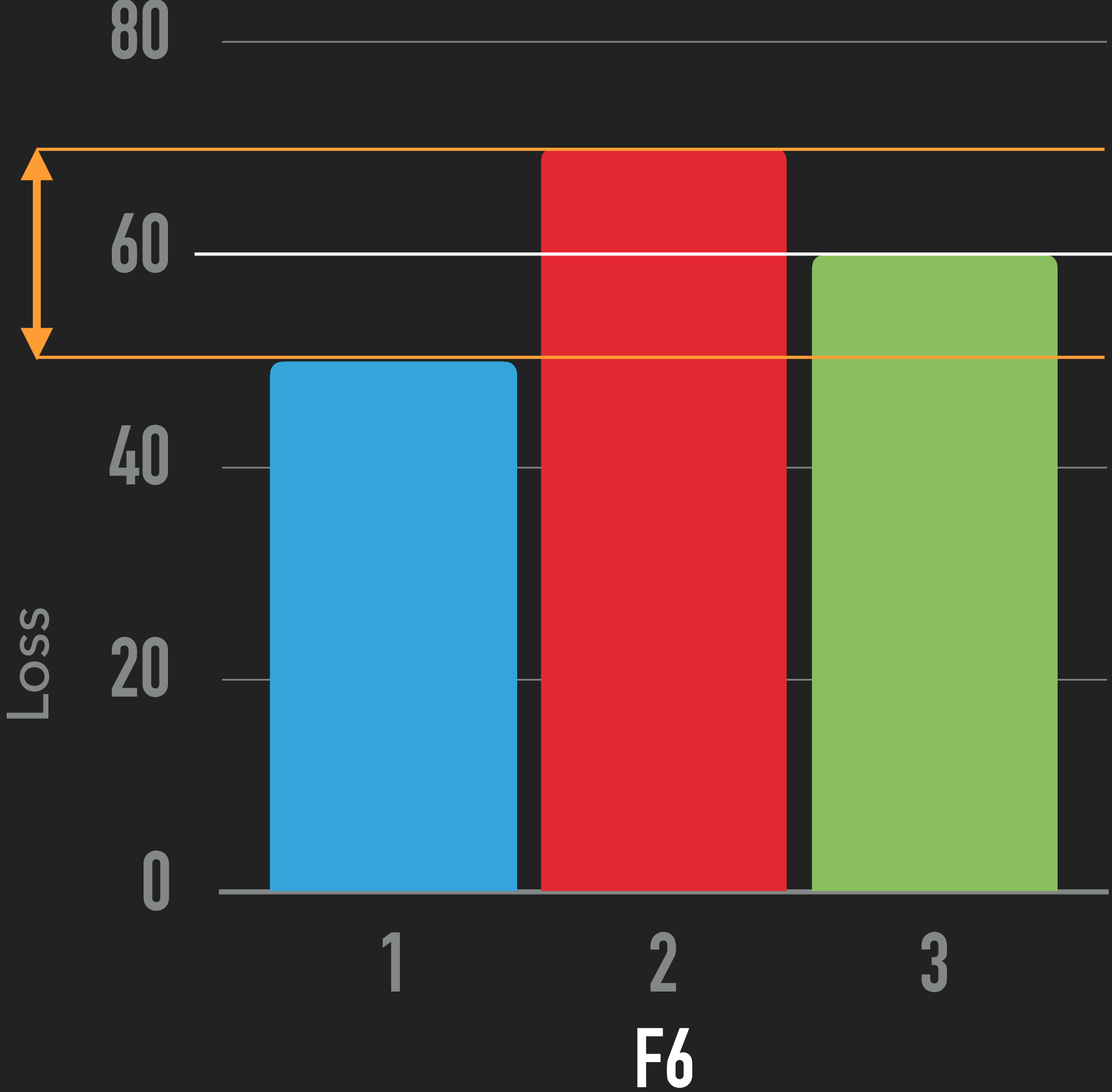
MINIMISING AGGREGATED EMPIRICAL RISK AND DEVIATION

▶ Trade off low deviation against higher average



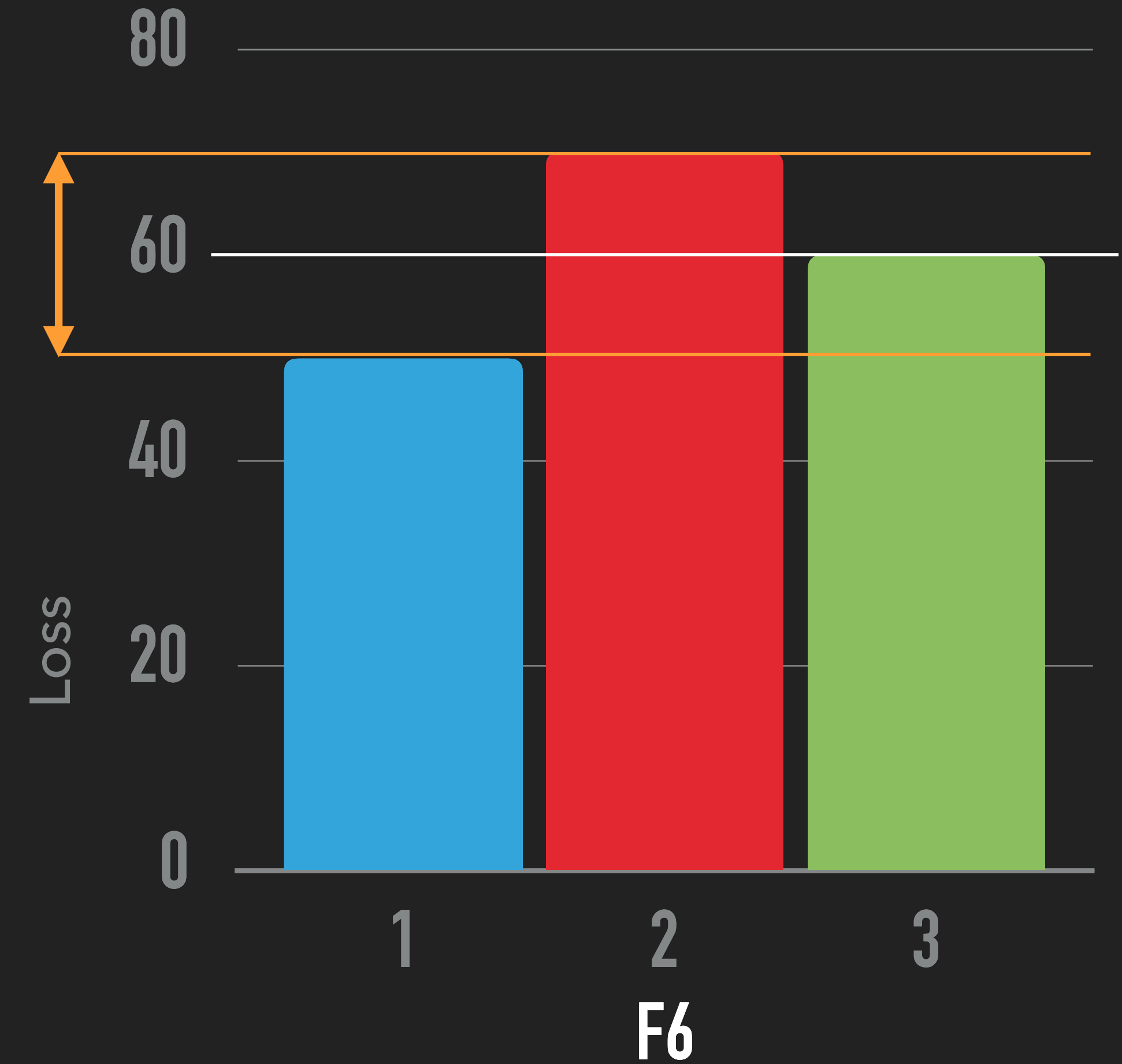
MINIMISING AGGREGATED EMPIRICAL RISK AND DEVIATION

▶ Trade off low **deviation** against higher **average**



MINIMISING AGGREGATED EMPIRICAL RISK AND DEVIATION

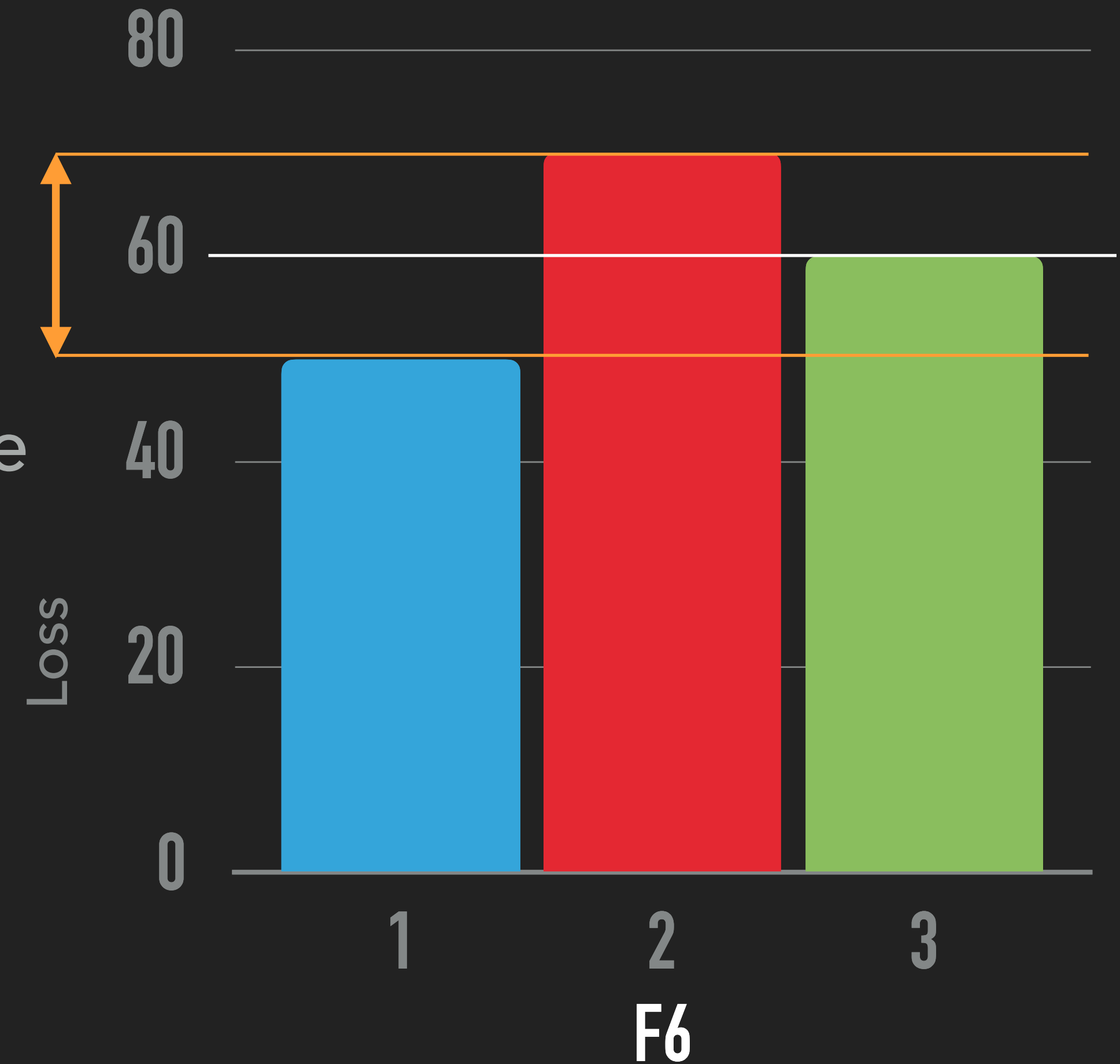
- ▶ Trade off low **deviation** against higher **average**
- ▶ Let $\mathcal{S} = \{1,2,3\}$ be the sensitive feature space



MINIMISING AGGREGATED EMPIRICAL RISK AND DEVIATION

- ▶ Trade off low **deviation** against higher **average**
- ▶ Let $\mathcal{S} = \{1,2,3\}$ be the sensitive feature space
- ▶ For $f \in \mathcal{F}$ let $R_f: \mathcal{S} \rightarrow \mathbb{R}$ be a r.v. (taking \mathcal{S} as the sample space, with a uniform base measure)

$$R_f: \mathcal{S} \ni s \mapsto \mathbb{E}_{(X,Y)} [\ell(Y, f(X)) \mid \mathcal{S} = s]$$



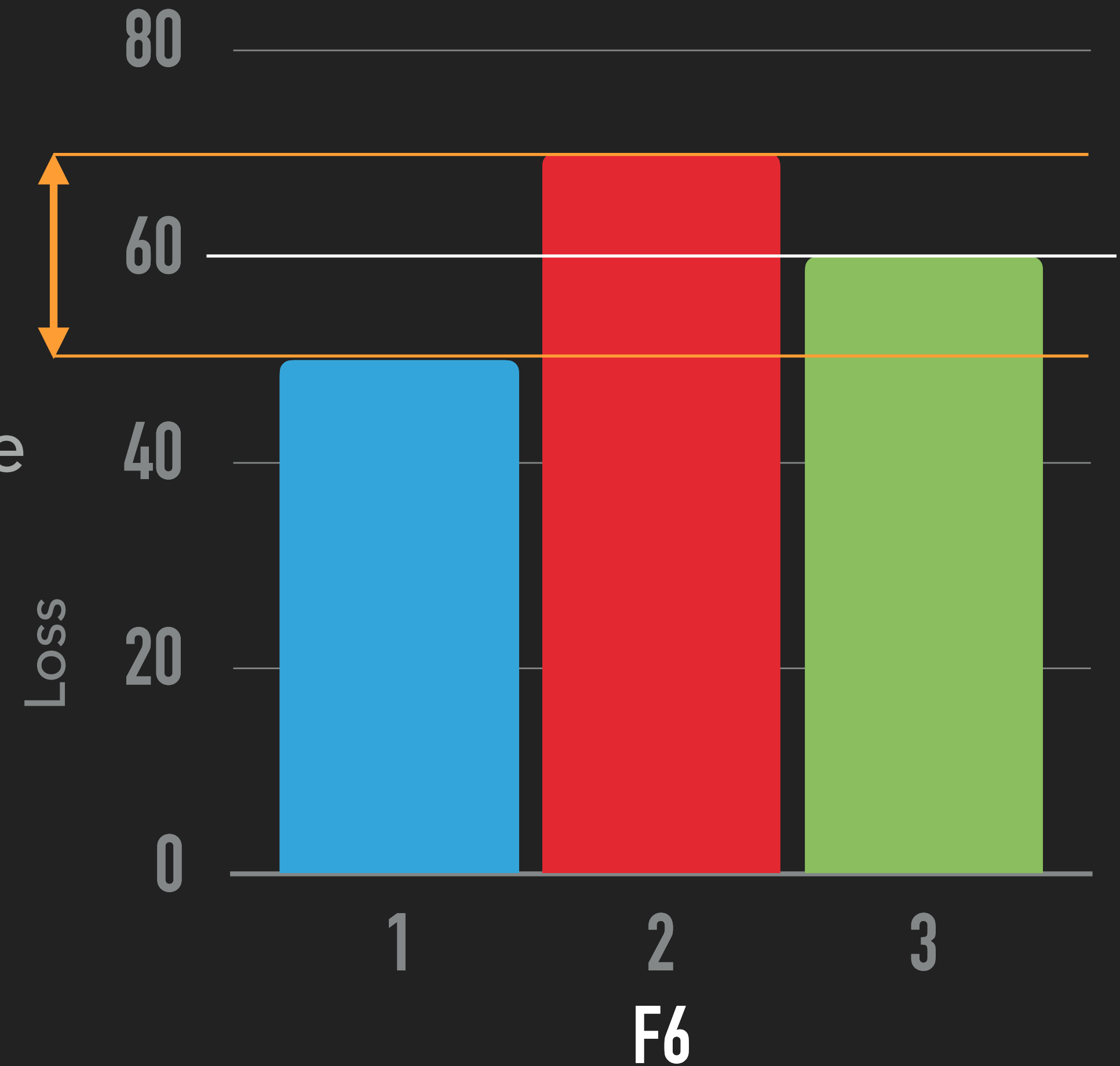
MINIMISING AGGREGATED EMPIRICAL RISK AND DEVIATION

- ▶ Trade off low **deviation** against higher **average**
- ▶ Let $\mathcal{S} = \{1,2,3\}$ be the sensitive feature space
- ▶ For $f \in \mathcal{F}$ let $R_f: \mathcal{S} \rightarrow \mathbb{R}$ be a r.v. (taking \mathcal{S} as the sample space, with a uniform base measure)

$$R_f: \mathcal{S} \ni s \mapsto \mathbb{E}_{(X,Y)} [\ell(Y, f(X)) \mid \mathcal{S} = s]$$

- ▶ Standard ERM:

$$\min_{f \in \mathcal{F}} \mathbb{E}(R_f)$$



MINIMISING AGGREGATED EMPIRICAL RISK AND DEVIATION

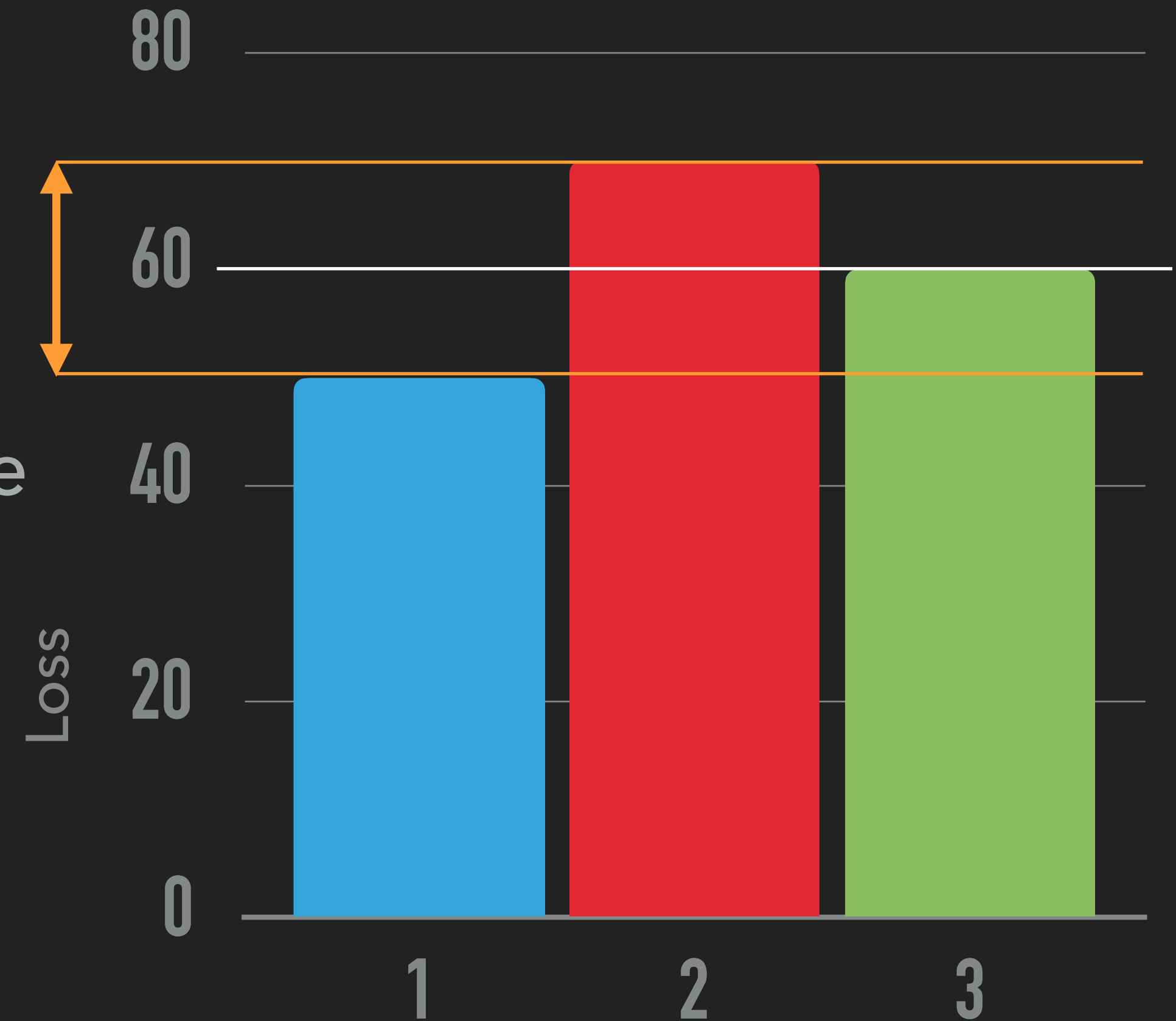
- ▶ Trade off low **deviation** against higher **average**
- ▶ Let $\mathcal{S} = \{1,2,3\}$ be the sensitive feature space
- ▶ For $f \in \mathcal{F}$ let $R_f: \mathcal{S} \rightarrow \mathbb{R}$ be a r.v. (taking \mathcal{S} as the sample space, with a uniform base measure)

$$R_f: \mathcal{S} \ni s \mapsto \mathbb{E}_{(X,Y)} [\ell(Y, f(X)) \mid \mathcal{S} = s]$$

- ▶ Standard ERM:

$$\min_{f \in \mathcal{F}} \mathbb{E}(R_f)$$

- ▶ Fairness Augmented ERM: $\min_{f \in \mathcal{F}} \mathbb{E}(R_f) + \mathcal{D}(R_f) = \min_{f \in \mathcal{F}} \mathcal{R}(R_f)$



F6

FAIRNESS RISK MEASURES ARE “REGULAR MEASURES OF RISK”

FAIRNESS RISK MEASURES ARE “REGULAR MEASURES OF RISK”

- ▶ Instead can **start** with axioms for \mathcal{R} Paper lists and justifies them

FAIRNESS RISK MEASURES ARE “REGULAR MEASURES OF RISK”

- ▶ Instead can **start** with axioms for \mathcal{R} Paper lists and justifies them
- ▶ Then show that such fairness risk measures are “regular measures of risk”
 - ▶ (In fact they are “coherent measures of risk”)

FAIRNESS RISK MEASURES ARE “REGULAR MEASURES OF RISK”

▶ Instead can start with axioms for \mathcal{R} Paper lists and justifies them

▶ Then show that such fairness risk measures are “regular measures of risk”

▶ (In fact they are “coherent measures of risk”)

▶ Such measures can always be written as

$$\mathcal{R}(R) = \mathbb{E}(R) + \mathcal{D}(R)$$

Fairness risk measure Deviation measure

FAIRNESS RISK MEASURES ARE “REGULAR MEASURES OF RISK”

▶ Instead can start with axioms for \mathcal{R} Paper lists and justifies them

▶ Then show that such fairness risk measures are “regular measures of risk”

▶ (In fact they are “coherent measures of risk”)

▶ Such measures can always be written as

$$\mathcal{R}(R) = \mathbb{E}(R) + \mathcal{D}(R)$$

Fairness risk measure Deviation measure

▶ Here \mathcal{D} is a “regular measure of deviation”

(i.e. convex, positively homogeneous, zero only when R is constant, and lower semicontinuous)

EXAMPLE RISK MEASURE AND CORRESPONDING DEVIATION MEASURE

EXAMPLE RISK MEASURE AND CORRESPONDING DEVIATION MEASURE

$$\mathcal{R}_{Q,\alpha}(Z) = \text{CVaR}_\alpha(Z) \quad \mathcal{D}_{Q,\alpha}(Z) = \text{CVaR}_\alpha(Z - \mathbb{E}(Z))$$

EXAMPLE RISK MEASURE AND CORRESPONDING DEVIATION MEASURE

$$\mathcal{R}_{Q,\alpha}(Z) = \text{CVaR}_\alpha(Z) \quad \mathcal{D}_{Q,\alpha}(Z) = \text{CVaR}_\alpha(Z - \mathbb{E}(Z))$$

- ▶ **CVaR** is the "Conditional Value at Risk".
- ▶ When Z is continuous random variable:

$$\text{CVaR}_\alpha(Z) = \mathbb{E}(Z | Z \geq q_\alpha(Z))$$

- ▶ where $q_\alpha(Z)$ is the α th quantile of Z

EXAMPLE RISK MEASURE AND CORRESPONDING DEVIATION MEASURE

$$\mathcal{R}_{Q,\alpha}(Z) = \text{CVaR}_\alpha(Z) \quad \mathcal{D}_{Q,\alpha}(Z) = \text{CVaR}_\alpha(Z - \mathbb{E}(Z))$$

- ▶ **CVaR** is the "Conditional Value at Risk".
- ▶ When Z is continuous random variable:

$$\text{CVaR}_\alpha(Z) = \mathbb{E}(Z \mid Z \geq q_\alpha(Z))$$

- ▶ where $q_\alpha(Z)$ is the α th quantile of Z
- ▶ Have $\text{CVaR}_0(Z) = \mathbb{E}(Z)$ and $\text{CVaR}_1(Z) = \max(Z)$

EXAMPLE RISK MEASURE AND CORRESPONDING DEVIATION MEASURE

$$\mathcal{R}_{Q,\alpha}(Z) = \text{CVaR}_\alpha(Z) \quad \mathcal{D}_{Q,\alpha}(Z) = \text{CVaR}_\alpha(Z - \mathbb{E}(Z))$$

- ▶ **CVaR** is the "Conditional Value at Risk".
- ▶ When Z is continuous random variable:

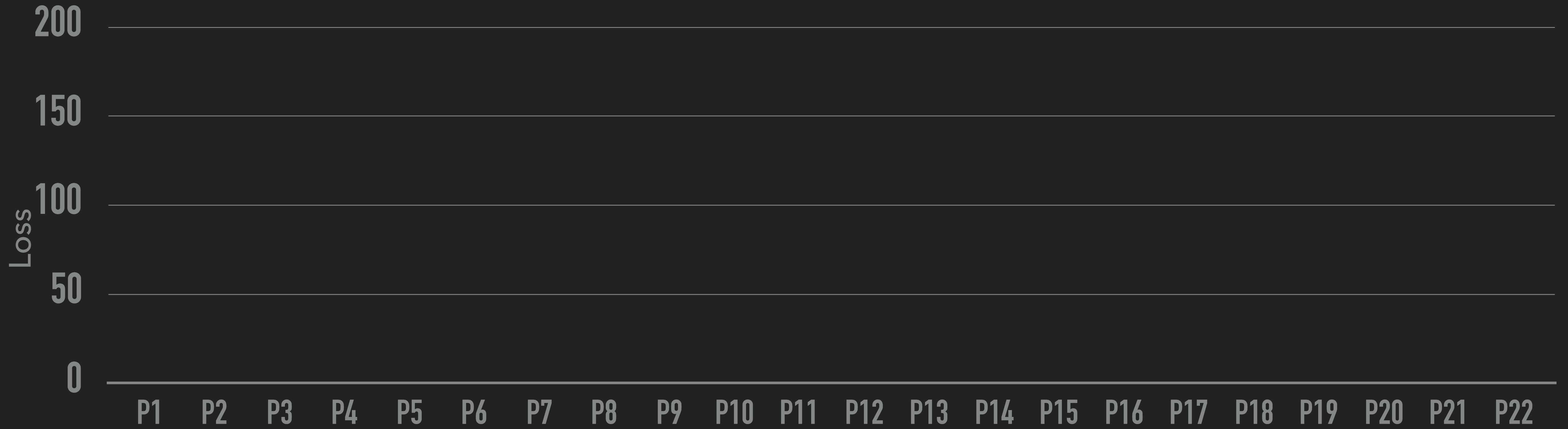
$$\text{CVaR}_\alpha(Z) = \mathbb{E}(Z \mid Z \geq q_\alpha(Z))$$

- ▶ where $q_\alpha(Z)$ is the α th quantile of Z
- ▶ Have $\text{CVaR}_0(Z) = \mathbb{E}(Z)$ and $\text{CVaR}_1(Z) = \max(Z)$
- ▶ Fairness objective becomes (see paper, eq (26)):

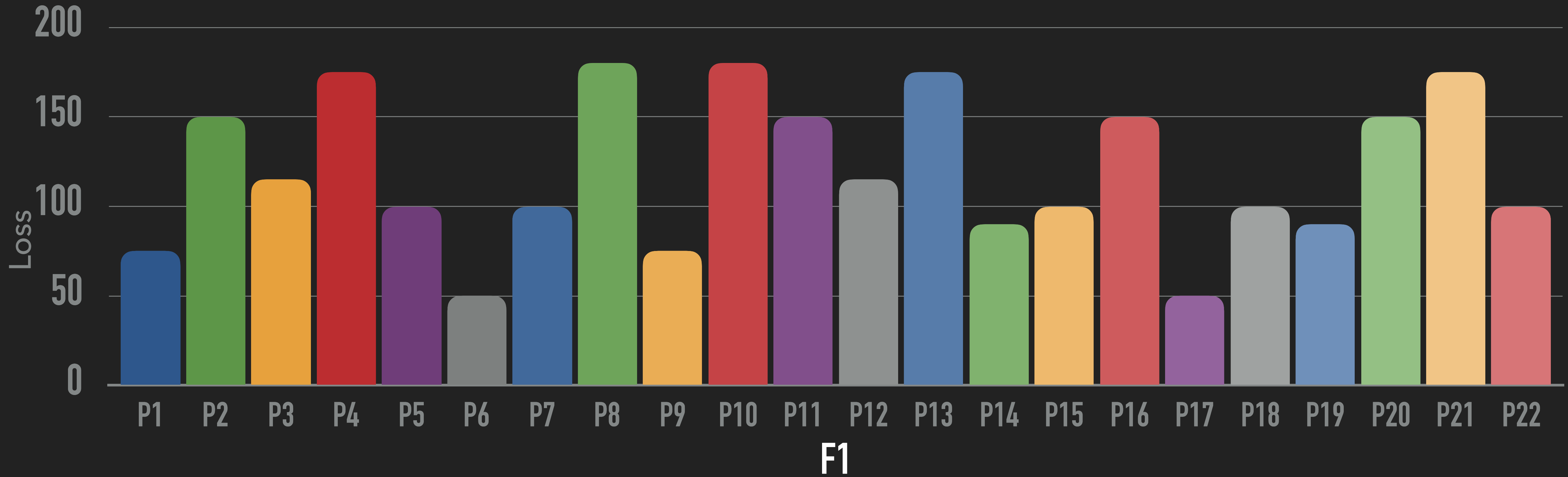
$$\min_{f \in \mathcal{F}, \rho \in \mathbb{R}} \left\{ \rho + \frac{1}{1-\alpha} \cdot \mathbb{E}[L(f) - \rho]_+ \right\}.$$

AN INTERESTING LIMITING CASE – EACH PERSON IS THEIR OWN CATEGORY!

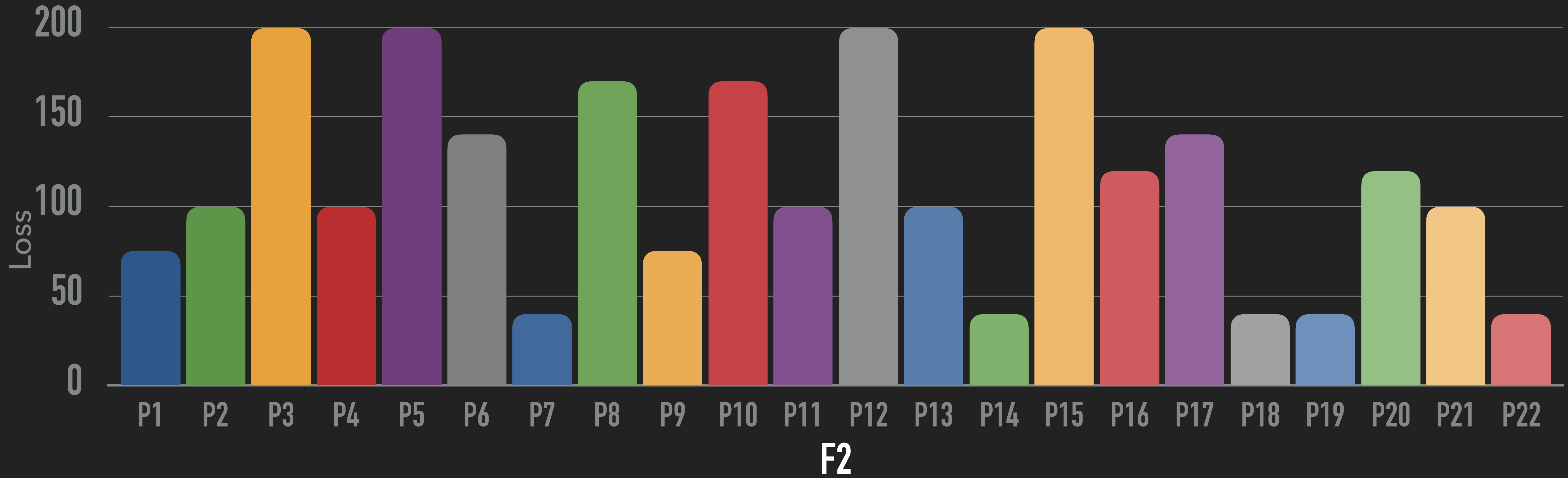
AN INTERESTING LIMITING CASE – EACH PERSON IS THEIR OWN CATEGORY!



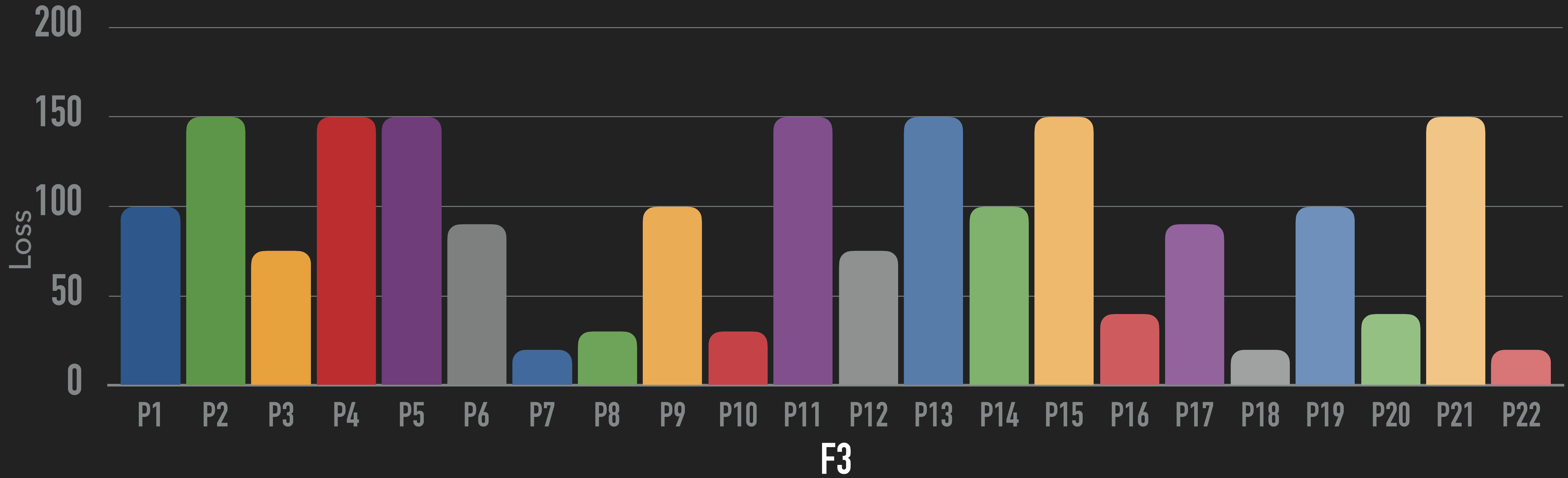
AN INTERESTING LIMITING CASE – EACH PERSON IS THEIR OWN CATEGORY!



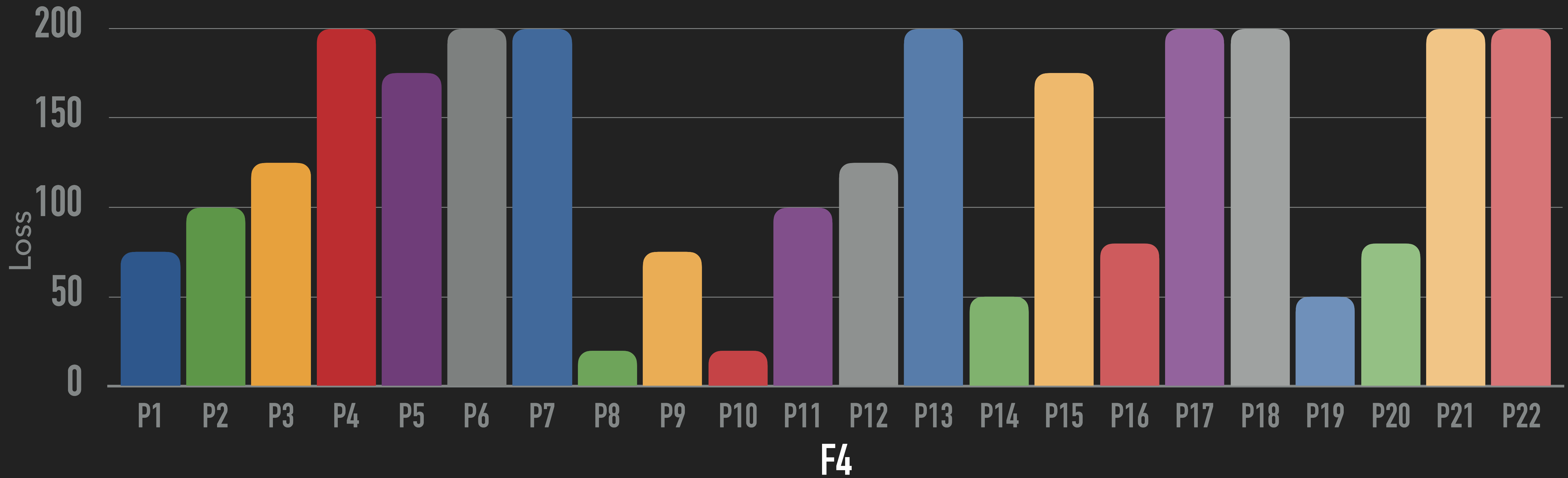
AN INTERESTING LIMITING CASE – EACH PERSON IS THEIR OWN CATEGORY!



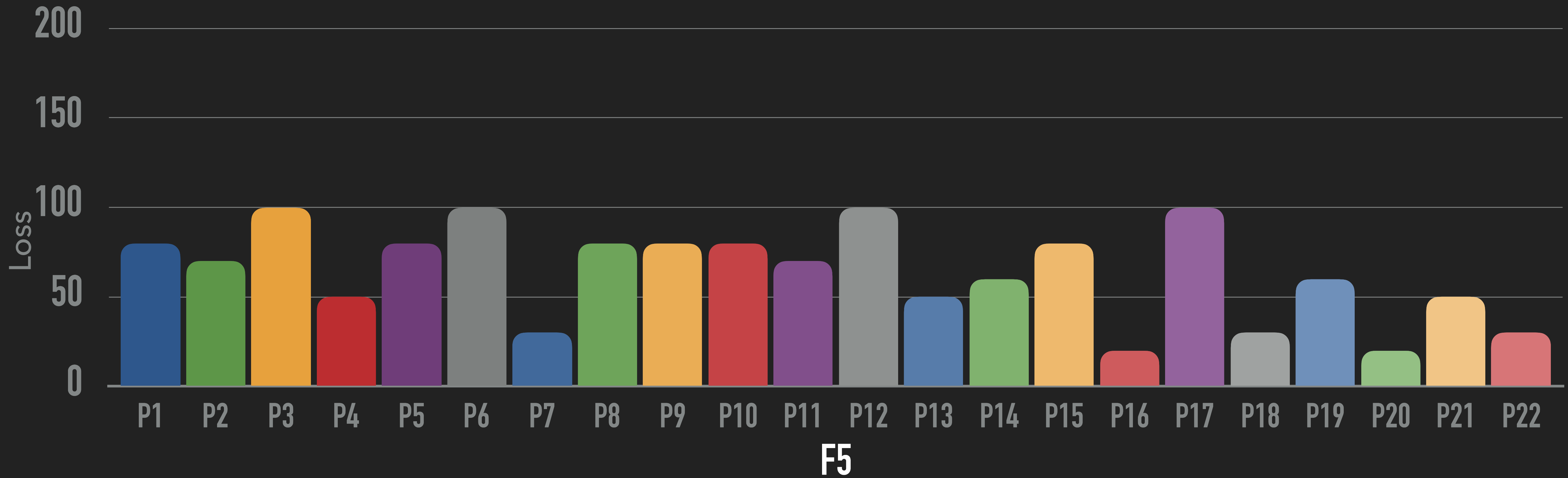
AN INTERESTING LIMITING CASE – EACH PERSON IS THEIR OWN CATEGORY!



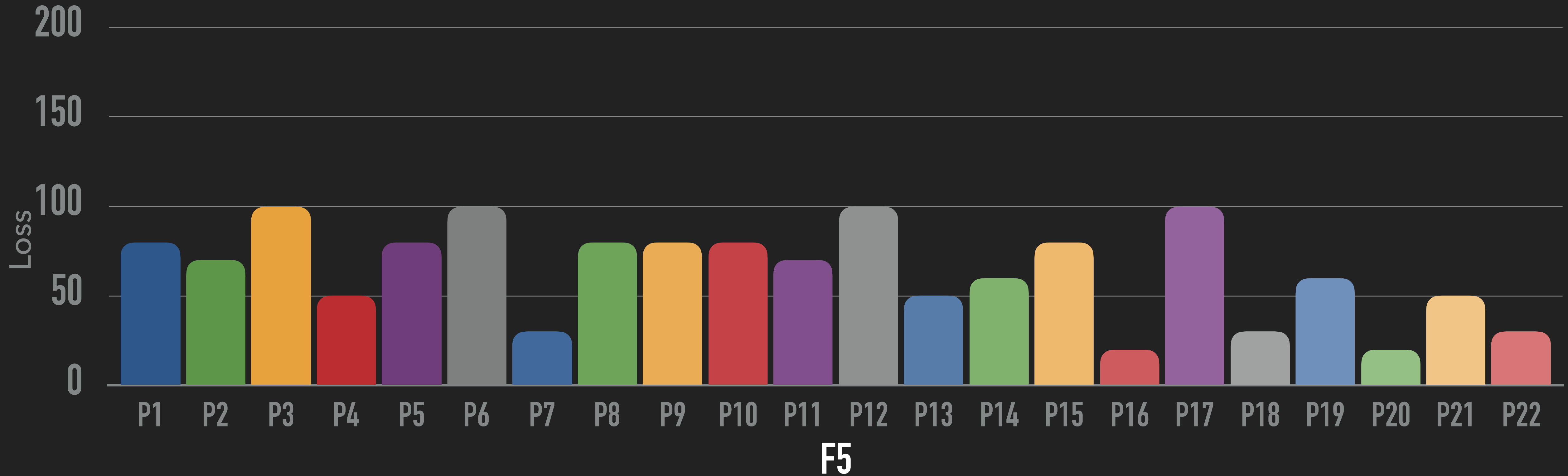
AN INTERESTING LIMITING CASE – EACH PERSON IS THEIR OWN CATEGORY!



AN INTERESTING LIMITING CASE – EACH PERSON IS THEIR OWN CATEGORY!

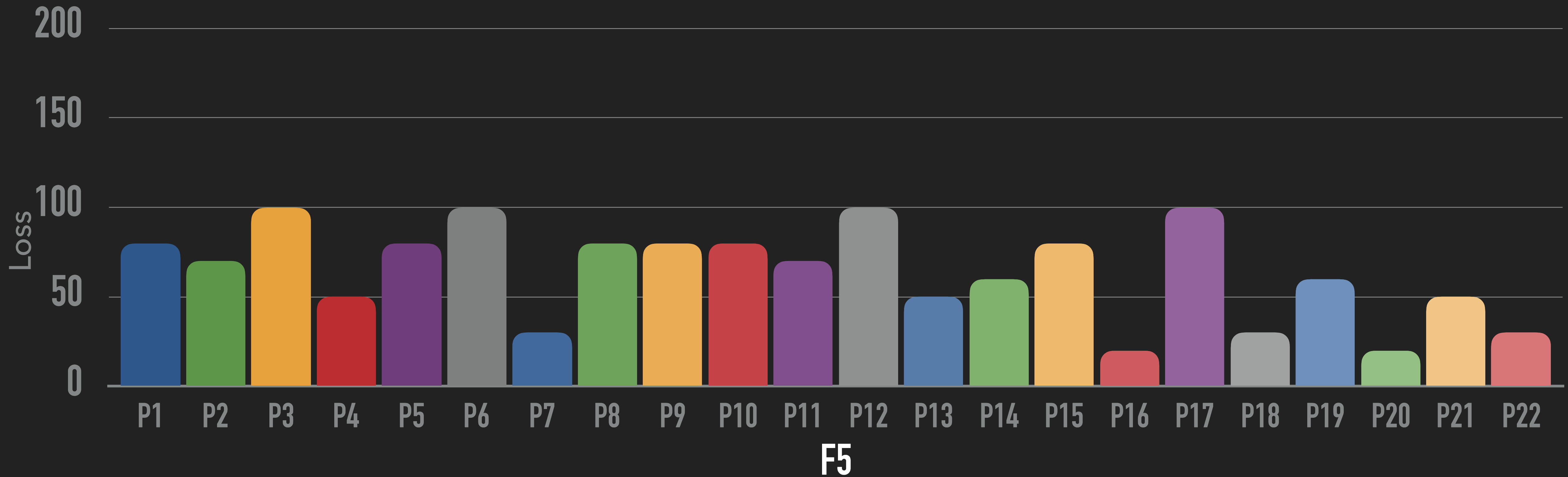


AN INTERESTING LIMITING CASE – EACH PERSON IS THEIR OWN CATEGORY!



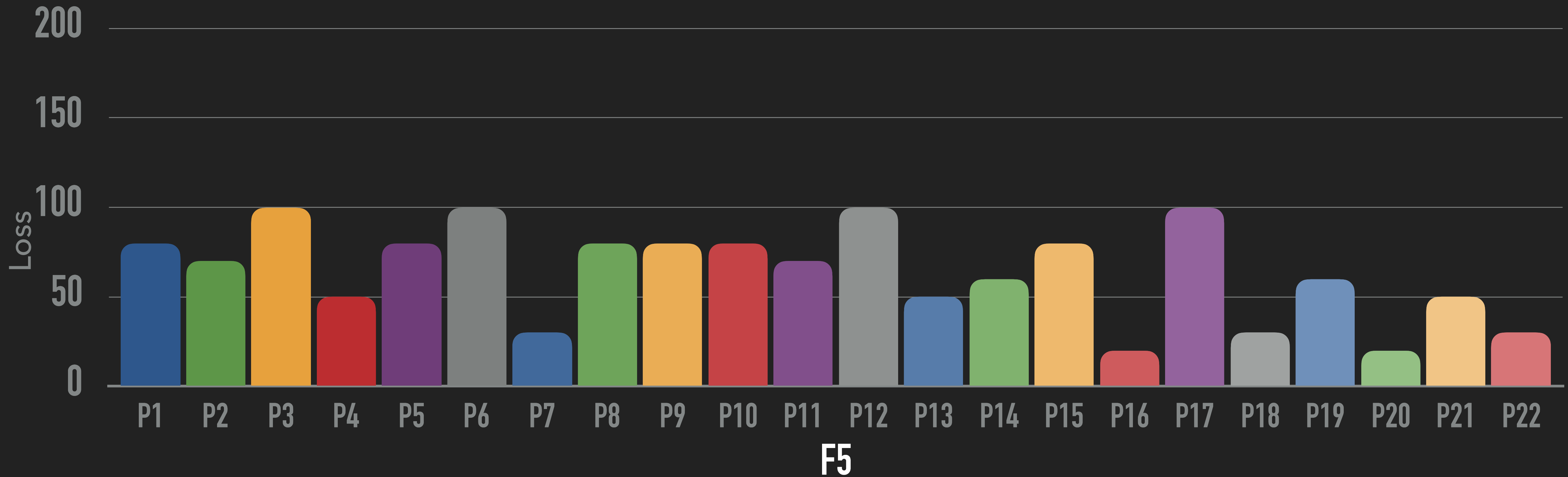
► Consistent with the principle that fundamental moral unit is the individual person

AN INTERESTING LIMITING CASE – EACH PERSON IS THEIR OWN CATEGORY!



- ▶ Consistent with the principle that fundamental moral unit is the individual person
- ▶ Avoids headaches with group boundaries and multiple group membership

AN INTERESTING LIMITING CASE – EACH PERSON IS THEIR OWN CATEGORY!



- ▶ Consistent with the principle that fundamental moral unit is the individual person
- ▶ Avoids headaches with group boundaries and multiple group membership
- ▶ Fairness risk measures automatically extend to this case (trivially)

CONCLUSION

CONCLUSION

$$\min_{f \in \mathcal{F}} \mathcal{R}(\mathbf{R}_f) \quad \mathbf{R}_f: \mathcal{S} \ni s \mapsto \mathbb{E}_{(X,Y)} [\ell(Y, f(X)) \mid \mathcal{S} = s]$$

- ▶ New and general approach to fairness in ML problems

CONCLUSION

$$\min_{f \in \mathcal{F}} \mathcal{R}(\mathbf{R}_f) \quad \mathbf{R}_f: \mathcal{S} \ni s \mapsto \mathbb{E}_{(X,Y)} [\ell(Y, f(X)) \mid \mathcal{S} = s]$$

- ▶ New and general approach to fairness in ML problems
- ▶ Fairness only depends upon **losses**, not **predictions**

CONCLUSION

$$\min_{f \in \mathcal{F}} \mathcal{R}(\mathbf{R}_f) \quad \mathbf{R}_f: \mathcal{S} \ni s \mapsto \mathbb{E}_{(X,Y)} [\ell(Y, f(X)) | \mathcal{S} = s]$$

- ▶ New and general approach to fairness in ML problems
- ▶ Fairness only depends upon **losses**, not **predictions**
- ▶ Fairness risk measures are **symmetric coherent measures of risk**

CONCLUSION

$$\min_{f \in \mathcal{F}} \mathcal{R}(\mathbf{R}_f) \quad \mathbf{R}_f: \mathcal{S} \ni s \mapsto \mathbb{E}_{(X,Y)} [\ell(Y, f(X)) | \mathcal{S} = s]$$

- ▶ New and general approach to fairness in ML problems
- ▶ Fairness only depends upon **losses**, not **predictions**
- ▶ Fairness risk measures are **symmetric coherent measures of risk**
- ▶ Close connection to **measures of inequality** (see appendix)

CONCLUSION

$$\min_{f \in \mathcal{F}} \mathcal{R}(\mathbf{R}_f) \quad \mathbf{R}_f: \mathcal{S} \ni s \mapsto \mathbb{E}_{(X,Y)} [\ell(Y, f(X)) \mid \mathcal{S} = s]$$

- ▶ New and general approach to fairness in ML problems
- ▶ Fairness only depends upon **losses**, not **predictions**
- ▶ Fairness risk measures are **symmetric coherent measures of risk**
- ▶ Close connection to **measures of inequality** (see appendix)
- ▶ Computationally **tractable**; related to **SVM!** (see paper / poster for **experiments**)

Humanising Machine Intelligence

Machine Learning Postdoc position available



hmi.anu.edu.au