# Risks with AI

- Overfitting

- Social and ethical risks

# Risks With AI: Overfitting

- Complex AI models such as Neural Nets can easily overfit (i.e. fit historical data too well but fail in realistic test conditions). If we don't understand what is helping the model perform well, there is a risk that the model will fail upon deployment

- Operational Risks
  - Direct financial risks (e.g. a trading algorithm)
  - Customer perception & reputation (e.g. poor personalization experience)

- ML Models need to go through thorough stress-testing (discussed later under audits)

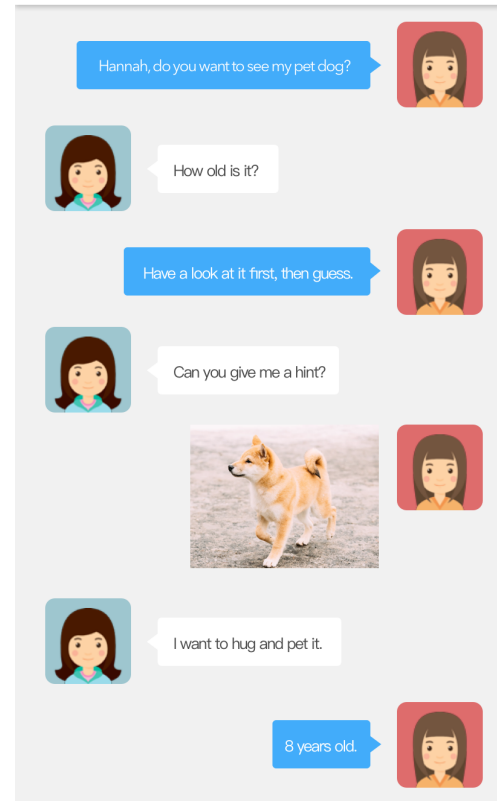# Xiaoice: Darling of Chinese Social Media



**XIAOBING**

40M followers

China



**YUAN ZHANG**

22 years old

China

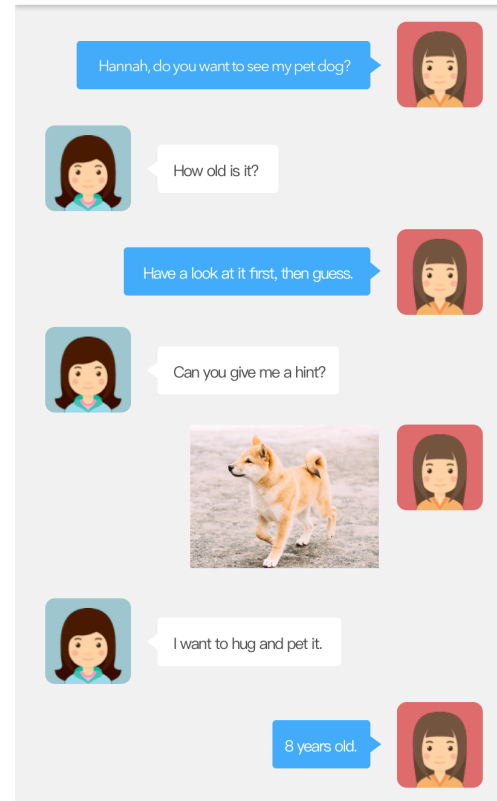# Xiaoice: Darling of Chinese Social Media



**XIAOICE**

40M followers

**Chatbot**

**YUAN ZHANG**

22 years old

China

# Tay.ai: Xiaoice's Evil Cousin



Microsoft's racist chatbot, Tay, makes MIT's annual worst-tech list
www.geekwire.com/2016/**microsoft**-chatbot-**tay**-mit-technology-fails/ ▾
Dec 27, 2016 - **Tay**, the **Microsoft** chatbot that pranksters trained to spew racist views, has resurfaced on MIT Technology Review's list of 2016's top technology ...
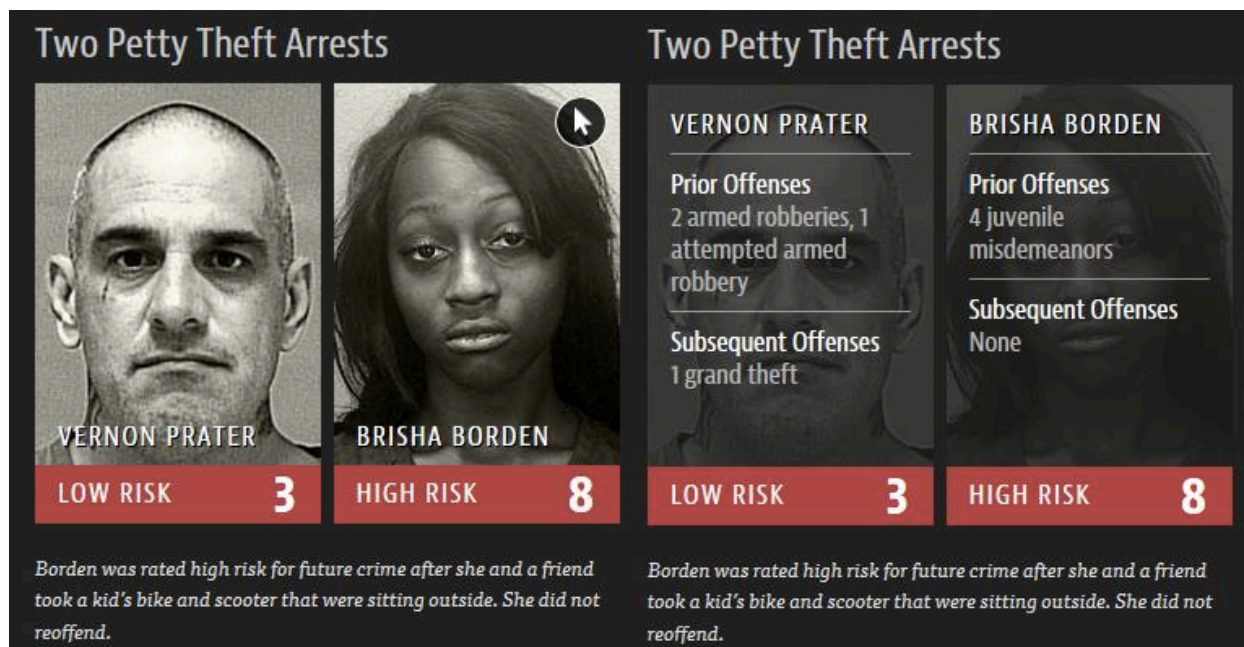
# Algorithm Bias in Recruiting Software



Amazon's machine learning specialists uncovered a big problem: their new recruiting engine **did not like women**.

# Algorithms Incorrectly Predict Recidivism

There are **higher false positive and false negative rates** among African Americans.
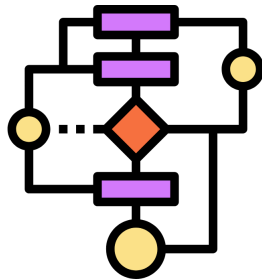


*Algorithms racial bias in predicting recidivism rates.*

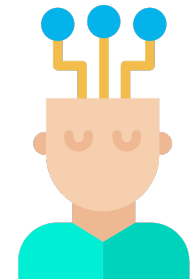# Why Might AI-based Decisions Unpredictable?

**Algorithm Logic**

**Nature**

**Data**

**Nurture**

**AI Behavior**



Analogous to human behavior, algorithms are driven by
**nature and nurture.**

# Risks to Society

- Social risks can result from automated decisions because these decisions may result in disadvantaged minorities continuing to be left behind.

  – Biases exist in the world and input data often contains these biases, so algorithms can incorporate discrimination and exacerbate it.

  – The harms that result can be categorized into the following two groups, originally described by Kate Crawford of the AI Now Institute.

| Harms of Allocation | Harms of Representation |
|---|---|
| • "When a system allocates or withholds opportunities or resources to or from certain groups"<br><br>• This affects what people can access<br><br>• E.g.: Unfairly denying access to a loan, mortgage, or job on the basis of gender, race, age, etc. | • "When a system misrepresents a group in a damaging way"<br><br>• This affects how people are viewed<br><br>• E.g.: Black faces labeled as gorillas by Google Photos; TSA scanners being more likely to false alarm people of color (based on hairstyles). |

# Risks to Firms

- These social risks then create additional risks for companies:

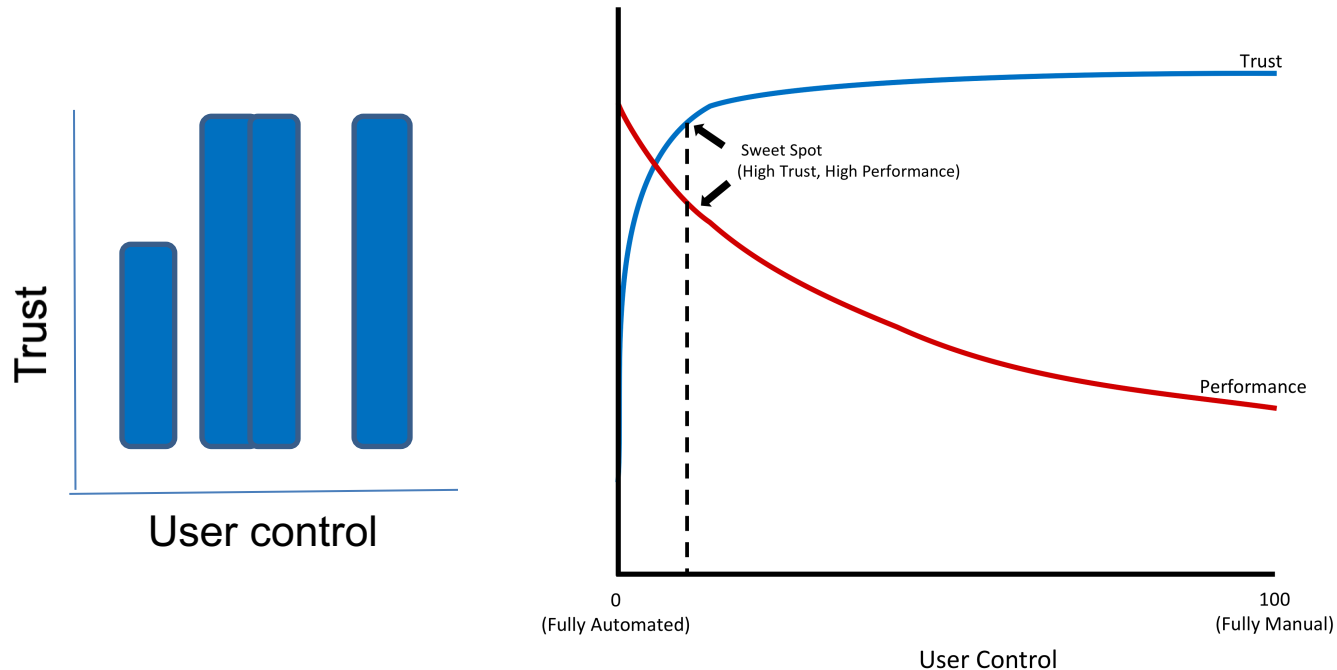| REPUTATIONAL RISK | LEGAL RISK | REGULATORY RISK |
|---|---|---|
| • Perceived to be a biased, prejudiced company<br><br>• Firms may face PR issues and backlash as a result | • Sued for unfair practices and discrimination against particular groups | • Increased regulation & cost of compliance<br><br>• Upcoming interest in auditing and data protection (GDPR) |

# AI Governance

- User Control

- Transparency

- Audits

# User Control

- Giving users some control over the decisions algorithms make for or about them
  - E.g., Facebook newsfeed now allows users to flag posts in their newsfeed as being false or offensive
  - System has helped detect many problematic posts

- In 2015, Facebook released mixed-style newsfeed controls
  - Allowed users to decide whether they wanted more/less of particular news (relationship statuses, profile changes, specific friends, etc.)
  - User satisfaction did increase (people liked having control)
  - But engagement went down (the algorithm knew what would engage users better than the users themselves)

- Demonstrates that there is a role for providing user control, but it needs to be balanced with ensuring the algorithm performs well.

Content/quotes from: "A Human's Guide to Machine Intelligence" by Kartik Hosanagar

# User Control

- There is a sweet spot that companies should aim for: providing a low level of control that enhances user trust, but not so much control that it significantly affects algorithm performance.



Content/quotes from: "A Human's Guide to Machine Intelligence" by Kartik Hosanagar
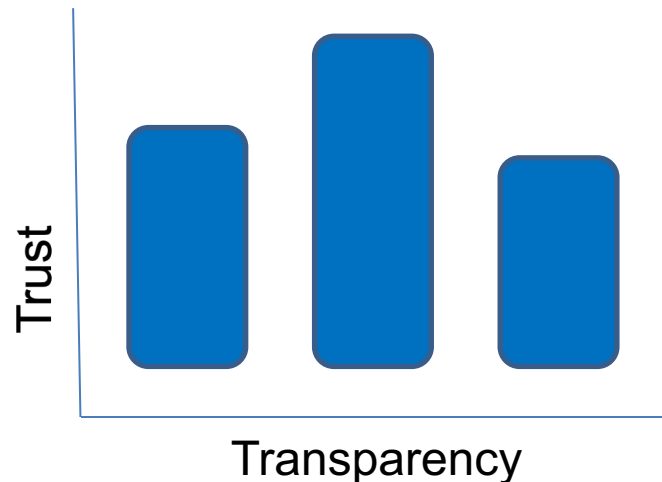
# Transparency for End Users

- Users can also be uncomfortable with decisions made by algorithms if they feel they don't understand those decisions.

- Solutions to this sometimes focus on **technical transparency**
  - Technical transparency = revealing source code
  - Examples: 2015 CFTC ruling on 2010 flash crash, 2017 NYC proposed bill on automated decisions
  - Problems with technical transparency: Doesn't protect intellectual property, creates vulnerability to adversarial attacks, doesn't always explain decisions (in cases of deep learning, random forests, etc.)

  - **Moreover, technical transparency may not even be necessary or helpful in improving user trust.**

**Wharton**
UNIVERSITY of PENNSYLVANIA

# Transparency for End Users

In a study on grading by Rene Kizilcec, students who were provided a limited amount of transparency had the greatest trust in their grades. Those who were given more extensive explanations of how the algorithm functioned had trust levels similar to those with no explanation at all.



**"Calibrated transparency"** is a better approach:

- Was an algorithm used to make a decision?

- What kinds of data are used?

- What variables are considered?

- Global and local interpretability

# Transparency for Managers and Data Scientists

- Many of the best performing ML models are often highly opaque (e.g. Neural Network making loan approval decisions)

- Recent interest in Interpretable ML (also known as Explainable AI)
  - Global interpretability: Can we explain at a high level what are the most important variables driving a model's predictions (e.g. income, credit history, etc)
  - Local interpretability: Can we explain the most important variables driving a particular prediction or decision (why was Kartik's loan application not approved?)

- Many open source tools and third-party vendors are starting to offer model interpretability solutions
  - Can be valuable for debugging

# Auditing Algorithms

- Due to the risks automated decisions bring, ML should be designated as a distinct model type, with it's own governance frameworks.

- The Algorithmic Accountability Act could be one such framework
  - If passed, it would require large companies to formally evaluate their "high-risk automated decision systems" for accuracy and fairness. But forward-thinking companies should not wait for regulation.

- What would an audit look like?
  - An audit process would begin with the creation of an inventory of all machine learning models being employed at a company as well as:
    - The specific uses of these models
    - The names of the developers & business owners of models
    - Risk ratings: the social/financial risks if the model fails

# Auditing Algorithms

- For high-risk models, the audit should look at the following three areas:

| Inputs | Model | Outputs |
|---|---|---|
| - Data quality<br>- Bias in training data | - Alternative models<br>- Statistical tests for model fit, overfitting<br>- Model transparency<br>- Stress test against simulated data | - Decisions with explanations<br>- Outliers: range of inputs and outputs |

- Machine learning models involve a variety of complex issues, such as bias, interpretability, and the fact that they are constantly retraining as they receive more data.

# Three Lines of Defense

Three Lines of Defense for Data Science



Model Developer

Data Science QA

Data science auditor
(for high-stakes models)